The Architecture of the Utrecht Haskell Compiler

Atze Dijkstra Jeroen Fokker S. Doaitse Swierstra

Department of Information and Computing Sciences
Universiteit Utrecht
P.O.Box 80.089, 3508 TB Utrecht, The Netherlands
{atze,jeroen,doaitse}@cs.uu.nl

Abstract

In this paper we describe the architecture of the Utrecht Haskell Compiler (UHC). UHC is a new Haskell compiler, that supports most (but not all) Haskell 98 features, plus some experimental extensions. It targets multiple backends, including a bytecode interpreter backend and a whole-program analysis backend, both via C. The implementation is rigorously organized as stepwise transformations through some explicit intermediate languages. The tree walks of all transformations are expressed as an algebra, with the aid of an Attribute Grammar based preprocessor. The compiler is just one materialization of a framework that supports experimentation with language variants, thanks to an aspect-oriented internal organization.

Categories and Subject Descriptors D.3.4 [Programming languages]: Compilers; Preprocessors; F.3.2 [Logics and meanings of programs]: Program analysis

General Terms Languages, Design

Keywords Haskell, compiler architecture, attribute grammar, aspect orientation

1. Introduction

On the occasion of the Haskell Hackathon on April 18th, 2009, we announced the first release of a new Haskell compiler: the Utrecht Haskell Compiler, or UHC for short. Until Haskell Prime [16] is available as a standard, UHC strives to be a full Haskell 98 [30] compiler (although currently it lacks a few features). The reason that we announce the compiler even though it is not yet fully finished, is that we feel that UHC is mature enough to use for play and experimentation.

One can ask why there is a need for (yet) another Haskell compiler, where the Glasgow Haskell Compiler (GHC) is already available as a widely used, fully featured, production quality Haskell compiler [26, 15, 28, 31]. In fact, we are using GHC ourselves for the implementation of UHC. Also, various alternatives exist, like Hugs (that in its incarnation of Gofer was the epoch maker for Haskell), and the Haskell compilers from York (NHC/YHC).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Haskell'09, September 3, 2009, Edinburgh, Scotland, UK. Copyright ⓒ 2009 ACM 978-1-60558-508-6/09/09...\$10.00

Still, we think UHC has something to add to existing compilers, not so much as a production compiler (yet), but more because of its systematically designed and extensible architecture. It is intended to be a platform for those who wish to experiment with adding new language or type system features. In a broader sense, UHC is a framework from which one can construct a series of increasingly complex compilers for languages reaching from simple lambda calculus to (almost-)Haskell 98. The UHC compiler in strict sense is just the culmination point of the series. We have been referring to the framework as 'EHC' (E for essential, extensible, educational, experimental...) in the past [10], but for ease we now call both the framework and its main compiler 'UHC'. Internally we use a stepwise and aspect-wise approach, realized by the use of attribute grammars (AG) and other tools.

In its current state, UHC supports most of the Haskell 98 (including polymorphic typing, type classes, input/output, base library), but a few features are still lacking (like defaulting, and some members of the awkward squad [29]). On the other hand, there are some extensions, notably to the type system. The deviations from the standard are not caused by obstinacy or desire to change the standard, but rather because of arbitrary priorization of the feature wish list.

The main structure of the compiler is shown in Figure 1. Haskell source text is translated to an executable program by stepwise transformation. Some transformations translate the program to a lower level language, many others are transformations within one language, establishing an invariant or performing an optimization.

All transformations, both within a language and between languages, are expressed as an algebra giving a semantics to the language. The algebras are described with the aid of an attribute grammar, which makes it possible to write multi-pass tree-traversals without even knowing the exact number of passes. Although the compiler driver is set up to pass data structures between transformations, for all intermediate languages we have a concrete syntax with a parser and a pretty printer. This facilitates debugging the compiler, by inspecting code between transformations.

Here is a short characterization of the intermediate languages. In section 3 we give a more detailed description.

- Haskell (HS): a general-purpose, higher-order, polymorphically typed, lazy functional language.
- Essential Haskell (EH): a higher-order, polymorphically typed, lazy functional language close to lambda-calculus, without syntactic sugar.
- Core: an untyped, lazy functional language close to lambdacalculus (at the time of this writing we are working on moving to a typed intermediate language, a combination of Henk [32], GHC core, and recent work on calling conventions [6]).

- Grin: 'Graph reduction intermediate notation', the instruction set of a virtual machine of a small functional language with strict semantics, with features that enable implementation of laziness [7].
- Silly: 'Simple imperative little language', an abstraction of features found in every imperative language (if-statements, assignments, explicit memory allocation) augmented with primitives for manipulating a stack, easily translatable to e.g. C (not all features of C are provided, only those that are needed for our purpose).
- BC: A bytecode language for a low-level machine intended to interpret Grin which is not whole-program analyzed nor transformed. We do not discuss this language in this paper.

The compiler targets different backends, based on a choice of the user. In all cases, the compiler starts compiling on a per module basis, desugaring the Haskell source text to Essential Haskell, type checking it and translating it to Core. Then there is a choice from three modes of operation:

- In whole-program analysis mode, the Core modules of the program and required libraries are assembled together and processed further as a whole. At the Grin level, elaborate intermodule optimization takes place. Ultimately, all functions are translated to low level C, which can be compiled by a standard compiler. As alternative backends, we are experimenting with other target languages, among which are the Common Intermediate Language (CIL) from the Common language infrastructure used by .NET [19], and the Low-Level Virtual Machine (LLVM) compiler infrastructure [25].
- In *bytecode interpreter mode*, the Core modules are translated to Grin separately. Each Grin module is translated into instructions for a custom bytecode machine. The bytecode is emitted in the form of C arrays, which are interpreted by a handwritten bytecode interpreter in C.
- In *Java mode*, the Core modules are translated to bytecode for the Java virtual machine (JVM). Each function is translated to a separate class with an *eval* function, and each closure is represented by an object combining a function with its parameters. Together with a driver function in Java which steers the interpretation, these can be stored in a Java archive (jar) and be interpreted by a standard Java interpreter.

The bytecode interpreter mode is intended for use during program development: it compiles fast, but because of the interpretation overhead the generated code is not very fast. The whole-program analysis mode is intended to use for the final program: it takes more time to compile, but generates code that is more efficient.

In Section 2 we describe the tools that play an important role in UHC: the Attribute Grammar preprocessor, a language for expressing type rules, and the variant and aspect manager. In Section 3 we describe the intermediate languages in the UHC pipeline in more detail, illustrated with a running example. In Section 4 the transformations are characterized in more detail. Finally, in Section 5 we draw conclusions about the methodology used, and mention related and future work.

2. Techniques and Tools

2.1 Tree-oriented programming

Using higher order functions on lists, like map, filter and foldr, is a good way to abstract from common patterns in functional programs.

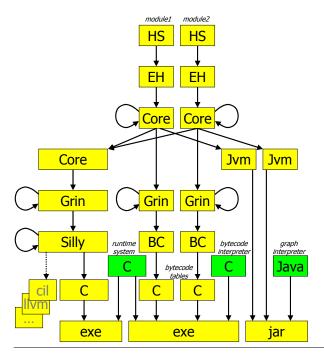


Figure 1. Intermediate languages and transformations in the UHC pipeline, in each of the three operation modes: whole-program analysis (left), bytecode interpreter (middle), and Java (right).

The idea that underlies the definition of foldr, i.e. to capture the pattern of an inductive definition by having a function parameter for each constructor of the data structure, can also be used for other data types, and even for multiple mutually recursive data types. A function that can be expressed in this way was called a catamorphism by Bird, and the collective extra parameters to foldr-like functions an algebra [3, 2]. Thus, ((+), 0) is an algebra for lists, and ((+), []) is another. In fact, every algebra defines a semantics of the data structure. When applying foldr-like functions to the algebra consisting of the original constructor functions, such as ((:), []) for lists, we have the identity function. Such an algebra is said to define the "initial" semantics. Outside circles of functional programmers and category theorists, an algebra is simply known as a "tree walk specification".

In compiler construction, algebras are very useful in defining a semantics of a syntactic structure or, bluntly said, to define tree walks over the parse tree. The fact that this is not widely done, is due to the following problems:

- 1. Unlike lists, for which *foldr* is standard, in a compiler we deal with custom data structures for abstract syntax of a language, which each need a custom *fold* function. Moreover, whenever we change the abstract syntax, we need to change the *fold* function and every algebra.
- 2. Generated code can be described as a semantics of the language, but often we need more than one alternative semantics: listings, messages, and internal structures (symbol tables etc.). This can be done by having the semantic functions in algebras return tuples, but this makes the program hard to maintain.
- Data structures for abstract syntax tend to have many alternatives, so algebras end up being clumsy tuples containing dozens of functions.
- 4. In practice, information not only flows bottom-up in the parse tree, but also top-down. E.g., symbol tables with global defini-

tions need to be distributed to the leaves of the parse tree to be able to evaluate them. This can be done by using higher-order domains for the algebras, but the resulting code becomes even harder to understand.

5. A major portion of the algebra is involved with moving information around. The essence of a semantics usually forms only a small part of the algebra and is obscured by lots of boilerplate.

Some seek the solution to these problems in the use of monads: the reader monad to pass information down into the tree, the writer monad to move information upwards, and the state monad and its derivatives to accumulate information during the tree walk [20]. Despite the attractiveness of staying inside Haskell we think this approach is doomed to fail when the algebras to be described are getting more and more complicated.

To save the nice idea of using an algebra for defining a semantics, we use a preprocessor [34] for Haskell that overcomes the above-mentioned problems. It is not a separate language; we can still use Haskell for writing auxiliary functions, and use all abstraction techniques and libraries available. The preprocessor just allows a few additional constructs, which can be translated into a custom *fold* function and algebras, or an equivalent more efficient implementation. (If one really wants to avoid a preprocessor, Viera, Swierstra and Swierstra recently described a technique to encode an attribute grammar directly in Haskell while keeping the advantages described below [35].)

We describe the main features of the preprocessor here, and explain why they overcome the five problems mentioned above. The abstract syntax of the language is defined in a **data** declaration, which is like an Haskell **data** declaration with named fields, however without the braces and commas. Constructor function names need not to be unique between types. As an example, consider a fragment of a typical imperative language:

```
data Stat
  = Assign \ dest :: String
                            src :: Expr
     While cond :: Expr
                             body :: Stat
     Group \ elems :: [Stat]
data Expr
  = Const num :: Int
     Var
            name :: String
     Add
            left :: Expr
                             right :: Expr
     Call
            name :: String
                            args :: [Expr]
```

The preprocessor generates corresponding Haskell data declarations (adding braces and commas, and making the constructors unique by prepending the type name, like $Expr_Const$), and generates a custom fold function. This overcomes problem 1 (except for the part that algebras change when sytax is changed, which will be solved below).

For any desired value we wish to compute over a tree, we can declare a "synthesized attribute". Possibly more than one data type can have the same attribute. For example, we can declare that both statements and expressions need to synthesize bytecode as well as listings, and that expressions can be evaluated to integer values:

```
attr Expr Stat syn bytecode :: [Instr] syn listing :: String attr Expr syn value :: Int
```

The preprocessor generates semantic functions that return tuples of synthesized attributes, but we can simply refer to attributes by name. This overcomes problem 2. Moreover, if at a later stage we add extra attributes, we do not have to refactor a lot of code.

The value of each attribute needs to be defined for every constructor of every data type which has the attribute. Such definitions are known as "semantic rules", and start with keyword sem.

```
 \begin{array}{lll} \mathbf{sem} \; Expr \mid Const \; \mathbf{lhs}.value = \; @num \\ \mid Add & \; \mathbf{lhs}.value = \; @left.value + \; @right.value \end{array}
```

This states that the synthesized (left hand side) value attribute of a Constant expression is just the contents of the num field, and that of an Add-expression can be computed by adding the value attributes of its subtrees. The @-symbol in this context should be read as "attribute", not to be confused with Haskell "as-patterns". At the left of the =-symbol, the attribute to be defined is mentioned; at the right, the defining Haskell expression is given. Each definition (or group of definitions) is labeled with a constructor (Const and Add in the example), which in turn are labeled with the datatype (Expr in the example). Vertical bars separate the constructors (and should not be confused with 'guarded' equations). The preprocessor collects and orders all definitions in a single algebra, replacing attribute references by suitable selections from the results of the tree walk on the children. This overcomes problem 3. To be able to pass information downward during a tree walk, we can

To be able to pass information downward during a tree walk, we can define "inherited" attributes (the terminology goes back to Knuth [22]). As an example, it can serve to pass down an environment, i.e. a lookup table that associates variables to values, which is needed to evaluate expressions:

```
\begin{split} & \textbf{type} \; Env = [(String, Int)] \\ & \textbf{attr} \; Expr \; \textbf{inh} \; env :: Env \\ & \textbf{sem} \; Expr \; | \; Var \; \textbf{lhs}.value = fromJust \; \$ \\ & \qquad \qquad lookup \; @\textbf{lhs}.env \; @name \end{split}
```

The preprocessor translates inherited attributes into extra parameters for the semantic functions in the algebra. This overcomes problem 4.

In many situations, sem rules only specify that attributes a tree node inherits should be passed unchanged to its children, as in a *Reader* monad. To scrap the boilerplate expressing this, the preprocessor has a convention that, unless stated otherwise, attributes with the same name are automatically copied. A similar automated copying is done for synthesized attributes passed up the tree, as in a *Writer* monad. When more than one child offers a synthesized attribute with the required name, we can specify to use an operator to combine several candidates:

```
attr Expr Stat syn listing use (#) []
```

which specifies that by default, the synthesized attribute *listing* is the concatenation of the *listing*s of all children that produce a sublisting, or the empty list if no child produces one. This overcomes problem 5, and the last bit of problem 1.

2.2 Rule-oriented programming

Using the attribute-grammar (AG) based preprocessor we can describe the part of a compiler related to tree walks concisely and efficiently. However, this does not give us any means of looking at such an implementation in a more formal setting. We use the domain specific language Ruler for describing the AG part related to the type system.

Although the use of Ruler currently is in flux because we are working on a newer version and therefore are only partially using Ruler for type system descriptions, we demonstrate some of its capabilities because it is our intent to tackle the difficulties involved with type system implementations by generating as much as possible automatically from higher level descriptions.

The idea of Ruler is to generate from a single source both a *LaTeX* rendering for human use in technical writing:

$$\begin{array}{c} v \text{ fresh} \\ \Gamma; \mathcal{C}^k; v \rightarrow \sigma^k \vdash^e e_1 : \sigma_a \rightarrow \sigma \rightsquigarrow \mathcal{C}_f \\ \Gamma; \mathcal{C}_f; \sigma_a \vdash^e e_2 : _ \rightsquigarrow \mathcal{C}_a \\ \hline \Gamma; \mathcal{C}^k; \sigma^k \vdash^e e_1 \ e_2 : \mathcal{C}_a \sigma \leadsto \mathcal{C}_a \end{array} \text{E.APP}_{HM}$$

and its corresponding AG implementation:

In this paper we neither further discuss the meaning or intention of the above fragments [9] nor explain Ruler [12] in depth. However, to sketch the underlying ideas we show the Ruler source code required for the above output; we need to define the scheme (or type) of a judgment and populate these with actual rules.

A scheme defines a LaTeX output template (judgeuse tex) with holes to be filled in by rules and a parsing template (judgespec).

```
 \begin{array}{l} \textbf{scheme} \ expr = \\ \textbf{holes} \ [\textbf{node} \ e : Expr, \textbf{inh} \ valGam : ValGam, \textbf{inh} \ knTy : Ty \\ \textbf{, thread} \ ty VarMp : \mathcal{C}, \textbf{syn} \ ty : Ty] \\ \textbf{judgeuse} \ \textbf{tex} \ valGam; ty VarMp. \textbf{inh}; knTy \\ \vdash ... "e" \ e : ty \leadsto ty VarMp. \textbf{syn} \\ \textbf{judgespec} \ valGam; ty VarMp. \textbf{inh}; knTy \\ \vdash e : ty \leadsto ty VarMp. \textbf{syn} \\ \end{array}
```

The rule for application is then specified by specifying premise judgments (judge above the dash) and a conclusion (below the dash) using the parsing template defined for scheme expr.

```
 \begin{array}{l} \textbf{rule} \ e.app = \\ \textbf{judge} \ tvarvFresh \\ \textbf{judge} \ expr = ty VarMp. \textbf{inh}; ty VarMp; (v \rightarrow knTy) \\ \qquad \vdash eFun : (ty.a \rightarrow ty) \leadsto ty VarMp. fun \\ \textbf{judge} \ expr = ty VarMp. fun; valGam; ty.a \\ \qquad \vdash eArg : ty.a \leadsto ty VarMp. arg \\ \qquad - \\ \textbf{judge} \ expr = ty VarMp. \textbf{inh}; valGam; knTy \\ \qquad \vdash (eFun \ eArg) \\ \qquad : (ty VarMp. arg \ ty) \leadsto ty VarMp. arg \\ \end{array}
```

For this example no further annotations are required to automatically produce AG code, except for the freshness of a type variable. The judgment tvarvFresh encapsulates this by providing the means to insert some handwritten AG code.

In summary, the basic idea of Ruler is to provide a description resembling the original type rule as much as possible, and then helping the system with annotations to allow the generation of an implementation and a LaTeX rendering.

2.3 Aspect-oriented programming

UHC's source code is organized into small fragments, each belonging to a particular *variant* and *aspect*. A variant represents a step in a sequence of languages, where each step adds some language features, starting with simply typed lambda calculus and ending with UHC. Each step builds on top of the previous one. Independent of a variant each step adds features in terms of aspects. For example, the

type system and code generation are defined as different aspects. UHC's build system allows for selectively building a compiler for a variant and a set of aspects.

Source code fragments assigned to a variant and aspects are stored in *chunked* text files. A tool called *Shuffle* then generates the actual source code when parameterized with the desired variant and aspects. Shuffle is language neutral, so all varieties of implementation languages can be stored in chunked format. For example, the following chunk defines a Haskell wrapper for variant 2 for the construction of a type variable:

```
%%[(2 hmtyinfer || hmtyast).mkTyVar
mkTyVar :: TyVarId -> Ty
mkTyVar tv = Ty_Var tv
%%]
```

The notation %[(2 hmtyinfer | hmtyast).mkTyVar begins a chunk for variant 2 with name mkTyVar for aspect hmtyinfer (Hindley-Milner type inference) or hmtyast (Hindley-Milner type abstract syntax), ended by %]. Processing by Shuffle then gives:

```
mkTyVar :: TyVarId \rightarrow Ty

mkTyVar \ tv = Ty\_Var \ tv
```

The subsequent variant 3 requires a more elaborate encoding of a type variable (we do not discuss this further). The wrapper must be redefined, which we achieve by explicitly overriding 2.mkTyVar by a chunk for 3.mkTyVar:

```
%%[(3 hmtyinfer || hmtyast).mkTyVar -2.mkTyVar
mkTyVar :: TyVarId -> Ty
mkTyVar tv = Ty_Var tv TyVarCateg_Plain
%%]
```

Although the type signature can be factored out, we refrain from doing so for small definitions.

Chunked sources are organized on a per file basis. Each chunked file for source code for UHC is processed by Shuffle to yield a corresponding file for further processing, depending on the language used. For chunked Haskell a single module is generated, for chunked AG the file may be combined with other AG files by the AG compiler.

The AG compiler itself also supports a notion of aspects, different from Shuffle's idea of aspects in that it allows definitions for attributes and abstract syntax to be defined independent of file and position in a file. Attribute definitions and attribute equations thus can be grouped according to the programmers sense of what should be together; the AG compiler combines all these definitions and generates corresponding Haskell code.

Finally, chunked files may be combined by Shuffle by means of explicit reference to the name of a chunk. This also gives a form of literate programming tools [23] where text is generated by explicitly combining smaller text chunks. For example, the above code for 2.mkTyVar and 3.mkTyVar are extracted from the chunked source code of UHC and combined with the text for this explanation by Shuffle.

3. Languages

The compiler translates a Haskell program to executable code by applying many small transformations. In the process, the program is represented using five different data structures, or languages. Some transformations map one of these languages to the next, some are transformations within one language. Together, the five languages span a spectrum from a full feature, lazy functional language, to a limited, low level imperative language.

3.1 The Haskell Language

The Haskell language (HS) closely follows Haskell's concrete syntax. A combinator-based, error-correcting parser parses the source text and generates an HS parse tree. It consists of numerous datatypes, some of which have many constructors. A *Module* consists of a name, exports, and declarations. Declarations can be varied: function bindings, pattern bindings, type signatures, data types, new types, type synonyms, class, instance... Function bindings involve a right hand side which is either an expression or a list of guarded expressions. An expression, in turn, has no less than 29 alternatives. All in all, the description of the context-free grammar consists of about 1000 lines of code.

We maintain sufficient information in the abstract syntax tree to reconstruct the original input, including layout and superfluous parentheses, with only the comments removed.

When processing HS we deal with the following tasks:

- Name resolution: Checking for properly introduced names and renaming all identifiers to the equivalent fully qualified names.
- Operator fixity and precedence: Expressions are parsed without taking into account the fixity and precedence of operators.
 Expressions are rewritten to remedy this.
- Name dependency: Definitions are reordered into different let bindings such that all identifier uses come after their definition. Mutually recursive definitions are put into one letrec binding.
- Definition gathering: Multiple definitions for the same identifier are merged into one.
- Desugaring: List comprehensions, do-notation, etc. are desugared.

In the remainder of this section on languages we use the following running example program to show how the various intermediate languages are used:

module M where

```
\begin{array}{l} len :: [a] \rightarrow Int \\ len \ [] = 0 \\ len \ (x : xs) = 1 + len \ xs \\ main = putStr \ (show \ (len \ (replicate \ 4 \ `x'))) \end{array}
```

3.2 The Essential Haskell Language

HS processing generates Essential Haskell (EH). The EH equivalent of the running example is shown below. Some details have been omitted and replaced by dots.

```
\begin{array}{l} \mathbf{let} \quad M.len :: [a] \to Int \\ \quad M.len \\ \quad = \lambda x_1 \to \mathbf{case} \ x_1 \ \mathbf{of} \\ \quad UHC.Prelude.[] \\ \quad \to UHC.Prelude.fromInteger \ 0 \\ \quad (UHC.Prelude.: x \ xs \ ) \\ \quad \to \dots \\ \mathbf{in} \\ \mathbf{let} \quad M.main = UHC.Prelude.putStr \dots \\ \mathbf{in} \\ \mathbf{let} \quad main :: UHC.Prelude.IO \dots \\ \quad main = UHC.Prelude.ehcRunMain \ M.main \\ \mathbf{in} \\ \mathbf{main} \end{array}
```

In constrast to the HS language, the EH language brings back the language to its essence, removing as much syntactic sugar as is

possible. An EH module consists of a single expression only, which is the body of the main function, with local let-bindings for the other top-level values.

Processing EH deals with the following tasks:

- *Type system:* Type analysis is done, types are erased when Core is generated. Type analysis can be done unhindered by syntactical sugar, error messages refer to the original source location but cannot reconstruct the original textual context anymore.
- Evaluation: Enforcing evaluation is made explicit by means of a let! Core construct.
- Recursion: Recursion is made explicit by means of a letrec Core construct.
- *Type classes:* All evidence for type class predicates are transformed to explicit dictionary parameters.
- Patterns: Patterns are transformed to their more basic equivalent, inspecting one constructor at a time, etc. .

3.3 The Core Language

The Core language is basically the same as lambda-calculus. The Core equivalent of the running example program is:

```
module M =
letrec
  \{M.len =
     \lambda M.x_1.1 \rightarrow
       let \{ 2 = M.x_1.1 \} in
        case \_2 of
          \{ C : \{...,...\} \rightarrow ...
          ; C[]\{\} \rightarrow
            let.
               \{ -3 =
                  (UHC.Prelude.packedStringToInteger)
                    (#String "0")} in
             let
               (UHC.Prelude.fromInteger)
                    (UHC.Prelude.\_d_1\_Num : DICT)
                    (_3)} in
             _4
in ...
```

A Core module, apart from its name, consists of nothing more than an expression, which can be thought of as the body of main:

```
data CModule
= Mod nm :: Name expr :: CExpr
```

An expression resembles an expression in lambda calculus. We have constants, variables, and lambda abstractions and applications of one argument:

Alternative Tup encodes a constructor, to be used with App to construct actual data alternatives or tuples. The Tag of a Tup encodes the Int tag, arity, and other information.

Furthermore, there is case distinction and local binding:

```
| Case \ expr :: CExpr \ alts :: [CAlt] \ dflt :: CExpr \ | Let \ categ :: Categ \ binds :: [CBind] \ body :: CExpr
```

The *categ* of a *Let* describes whether the binding is recursive, strict, or plain. These two constructs use the auxiliary notions of alternative and binding:

```
 \begin{array}{ll} \mathbf{data} \ CAlt \\ &= Alt \quad pat : CPat \quad expr :: CExpr \\ \mathbf{data} \ CBind \\ &= Bind \ name : Name \ expr :: CExpr \\ &\mid \ FFI \quad name : Name \ imp :: String \ ty :: Ty \end{array}
```

A pattern introduces bindings, either directly or as a field of a constructor:

```
 \begin{array}{lll} \textbf{data} & CPat \\ &= Var & name :: Name \\ &| Con & name :: Name \ tag :: Tag \ binds :: [CPatBind] \\ &| BoolExpr \ name :: Name \ cexpr :: CExpr \\ \textbf{data} & CPatBind \\ &= Bind & offset :: Int & pat :: CPat \end{array}
```

The actual Core language is more complex because of:

- Experiments with extensible records; we omit this part as extensible records are currently not supported in UHC.
- Core generation is partly non syntax directed because context reduction determines which dictionaries are to be used for class predicates. The syntax directed part of Core generation therefore leaves holes, later to be filled in with the results of context reduction; this is a mechanism similar to type variables representing yet unknown types.
- An annotation mechanism is used to propagate information about dictionary values. This mechanism is somewhat ad hoc and we expect it to be changed when more analyses are done in earlier stages of the compiler.

3.4 The Grin Language

The Grin equivalent of the running example program is:

```
 \begin{array}{l} \textbf{module } M \\ \{\textit{M.len } \textit{M.x}_1.1 = \\ \{\textbf{eval } \textit{M.x}_1.1; \lambda.2 \rightarrow \\ \textbf{case .2 of} \\ \{\textit{C} \ / : \rightarrow \\ \{...\} \\ ; \textit{C} \ / \ [] \rightarrow \\ \{\textbf{store } (\textit{C} \ / \textit{UHC.Prelude.PackedString "0"}); \lambda.6 \rightarrow \\ \textbf{store } (\textit{F} \ / \textit{UHC.Prelude.packedStringToInteger .6}); \\ \lambda.3 \rightarrow \\ \textbf{store } (\textit{P} \ / \textit{O} \ / \textit{UHC.Prelude.fromInteger} \\ \textit{UHC.Prelude..d}_1.\textit{Num}); \lambda.5 \rightarrow \\ \textbf{store } (\textit{A} \ / \textit{apply .5 .3}); \lambda.4 \rightarrow \\ \textbf{eval } A \ \} \\ \}\} \} \end{array}
```

A Grin module consists of its name, global variables with their initializations, and bindings of function names with parameters to their bodies.

```
 \begin{array}{l} \textbf{data} \ GrModule \\ = \ Mod \ nm :: Name \ globals :: [\ GrGlobal] \ binds :: [\ GrBind] \\ \textbf{data} \ GrGlobal \\ = \ Glob \ nm :: Name \ val :: GrVal \end{array}
```

```
data GrBind
```

```
= \textit{Bind nm} :: \textit{Name args} :: [\textit{Name}] \textit{ body} :: \textit{GrExpr}
```

Values manipulated in the Grin language are varied: we have nodes (think: heap records) consisting of a tag and a list of fields, standalone tags, literal ints and strings, pointers to nodes, and 'empty'. Some of these are directly representable in the languages (nodes, tags, literal ints and strings)

```
 \begin{aligned} \textbf{data} & \ GrVal \\ &= \ LitInt & \ int :: Int \\ & \ \mid \ LitStr & \ str :: String \\ & \ \mid \ Tag & \ tag :: GrTag \\ & \ \mid \ Node & \ tag :: GrTag & \ flds :: [\ GrVal\ ] \end{aligned}
```

Pointers to nodes are also values, but they have no direct denotation. On the other hand, variables ranging over values are not a value themselves, bur for syntactical convenience we do add the notion of a 'variable' to the GrVal data type:

```
| Var name :: Name
```

The tag of a node describes its role. It can be a constructor of a datatype (Con), a function of which the call is deferred because of lazy evaluation (Fun), a function that is partially applied but still needs more arguments (PApp), or a deferred application of an unknown function (appearing as the first field of the node) to a list of arguments (App).

```
 \begin{array}{lll} \textbf{data} \ GrTag \\ &= \ Con & name :: Name \\ &| \ Fun & name :: Name \\ &| \ PApp \ needs :: Int \ name :: Name \\ &| \ App & apply fn :: Name \end{array}
```

The four tag types are represented as C, F, P and A in the example program above.

The body of a function denotes the calculation of a value, which is represented in a program by an 'expression'. Expressions can be combined in a monadic style. Thus we have *Unit* for describing a computation immediately returning a value, and *Seq* for binding a computation to a variable (or rather a lambda pattern), to be used subsequently in another computation:

```
 \begin{array}{ll} \mathbf{data} \ \mathit{GrExpr} \\ = \ \mathit{Unit} \ \mathit{val} :: \mathit{GrVal} \\ \mid \ \mathit{Seq} \ \ \mathit{expr} :: \mathit{GrExpr} \ \mathit{pat} :: \mathit{GrPatLam} \ \mathit{body} :: \mathit{GrExpr} \end{array}
```

There are some primitive computations (that is, constants in the monad) one for storing a node value (returning a pointer value), and two for fetching a node previously stored, and for fetching one field thereof:

```
 | Store val :: GrVal \\ | FetchNode name :: Name \\ | FetchField name :: Name offset :: Int
```

Other primitive computations call Grin and foreign functions, respectively. The name mentioned is that of a known function (i.e., there are no function variables) and the argument list should fully saturate it:

```
| Call \ name :: Name \ args :: [GrVal] | FFI \ name :: String \ args :: [GrVal]
```

Two special primitive computations are provided for evaluating node that may contain a Fun tag, and for applying a node that must contain a PApp tag (a partially applied function) to further arguments:

```
 \mid Eval \ name :: Name \\ \mid App \ name :: Name \ args :: [\ GrVal\ ]
```

Next, there is a computation for selecting a matching alternative, given the name of the variabele containing a node pointer:

```
| Case \ val :: GrVal \ alts :: [GrAlt]
```

Finally, we need a primitive computation to express the need of 'updating' a variable after it is evaluated. Boquist proposed an *Update* expression for the purpose which has a side effect only and an 'empty' result value [7]. We observed that the need for updates is always next to either a *FetchNode* or a *Unit*, and found it more practical and more efficient to introduce two update primitives:

```
| FetchUpdate \ src :: Name \ dst :: Name \ UpdateUnit \ name :: Name \ val :: GrVal
```

Auxiliary data structures are that for describing a single alternative in a *Case* expression:

and for two kinds of patterns, occurring in a Seq expression and in an Alt alternative, respectively. A simplified version of these is the following, but in reality we have more pattern forms.

```
 \begin{array}{ll} \mathbf{data} \ GrPatLam \\ = \ Var \quad name :: Name \\ \mathbf{data} \ GrPatAlt \\ = \ Node \ tag :: GrTag \ args :: [Name] \end{array}
```

4. Transformations

An UHC architecture principle is that the program is transformed in many small steps, each performing an isolated task. Even when multiple steps could have been combined, we prefer the simplicity of doing one task at a time. The Attribute Grammar preprocessor makes the definition of a tree walk easy, and the runtime overhead for the additional passes is modest.

Currently we have 12 transformations on the Core language, 24 on the Grin language, and 4 on the Silly language. Some of them are applied more than once, so the total number of transformations a program undergoes is even larger. In this section we give a short description of all transformations. Of course, this is just a snapshot of the current situation: the very fact that the steps are isolated and identified enables us to move them around while developing the compiler. Yet, the description of the transformations gives an idea of the granularity of the steps, and as a whole gives an overview of techniques employed.

4.1 Core Transformations

Three major gaps have to be bridged in the transformation from Core to Grin. Firstly, where Core has a lazy semantics, in Grin deferring of function calls and their later evaluation is explicitly encoded. Secondly, in Core we can have local function definitions, whereas in Grin all function definitions are at top level. Grin does have a mechanism for local, explicitly sequenced variable bindings. Thirdly, whereas Core functions always have one argument, in Grin functions can have multiple parameters, but they take them all at the same time. Therefore a mechanism for partial parametrization is necessary. The end result is lambda lifted Core, that is the floating of lambda-expressions to the top level and passing of non-global variables explicitly as parameters.

Core has one construct let! for enforcing evaluation to WHNF independent of other Core language constructs. This makes the implementation of seq easier but burdens Core transformations with the need not to cross an 'evaluation boundary' when moving code around.

The Core transformations listed below also perform some trivial cleanup and optimizations, because we avoid burdening the Core generation from EH with such aspects.

- 1. *EtaReduction* Performs restricted η -reduction, that is replace expressions like $\lambda x \ y \to f \ x \ y$ with f with the restriction that arity is not changed. Such expressions are introduced by coercions which (after context reduction) turn out not to coerce anything at all.
- 2. RenameUnique Renames variables such that all variables are globally unique.
- 3. LetUnrec Replaces mutually recursive bindings

letrec{
$$v_1 = \ldots; v_2 = \ldots$$
} in..

which actually are not mutually recursive by plain bindings

let
$$v_1 = \dots$$
 in let $v_2 = \dots$ in \dots

Such bindings are introduced because some bindings are potentially mutually recursive, in particular groups of dictionaries.

- 4. InlineLetAlias Inlines let bindings for variables and constants.
- 5. ElimTrivApp Eliminates application of the id function.
- 6. ConstProp Performs addition of int constants at compile time.
- 7. ANormal Complex expressions like

are broken up into a sequence of bindings and simpler expressions

let
$$v_1 = g$$
 a in let $v_2 = h$ b in f v_1 v_2

which only have variable references as their subexpressions.

- 8. *LamGlobalAsArg* Pass global variables of let-bound lambdaexpressions as explicit parameters, as a preparation for lambdalifting.
- CAFGlobalAsArg Similar for let-bound constant applicative forms (CAFs).
- 10. *FloatToGlobal* Performs 'lambda lifting': move bindings of lambda-expressions and CAFs to the global level.
- 11. *LiftDictFields* Makes sure that all dictionary fields exist as a top-level binding.
- 12. *FindNullaries* Finds nullary (parameterless) functions f and inserts another definition f' = f, where f' is annotated in such a way that it will end up as an updateable global variable.

After the transformations, translation to Grin is performed, where the following issues are addressed:

- for Let-expressions: global expressions are collected and made into Grin function bindings; local non-recursive expressions are sequenced by Grin Seq-expressions; for local recursive letbindings a Sequence is created which starts out to bind a new variable to a 'black hole' node, then processes the body, and finally generates a FetchUpdate-expression for the introduced variable.
- for Case-expressions: an explicit Eval-expression for the scrutinee is generated, in Sequence with a Grin Case-expression.
- for App-expressions: it is determined what it is that is applied:
 - if it is a constructor, then a node with *Con* tag is returned;
 - if it is a lambda of known arity which has exactly the right number of arguments, then either a Call-expression

- is generated (in strict contexts) or a node with Fun tag is stored with a Store-expression (in lazy contexts);
- if it is a lambda of known arity that is undersaturated (has not enough arguments), then a node with PApp tag is returned (in strict contexts) or Stored (in lazy contexts)
- if it is a lambda of known arity that is oversaturated (has too many arguments), then (in strict contexts) first a *Call*-expression to the function is generated that applies the function to some of the arguments, and the result is bound to a variable that is sub*Seq*uently *Applied* to the remaining arguments; or (in non-strict contexts) a node with *Fun* tag is *Stored*, and bound to a variable that is used in another node which has an *App* tag.
- if it is a variable that represents a function of unknown arity, then (in strict contexts) the variable is explicitly *Eval*uated, and its result used in an *App* expression to the arguments; or (in non-strict contexts) as a last resort, both function variable and arguments are stored in a node with *App* tag.
- for global bindings: lambda abstractions are 'peeled off' the body, to become the arguments of a Grin function binding.
- for foreign function bindings: functions with *IO* result type are treated specially.

We have now reached the point in the compilation pipeline where we perform our whole-program analysis. The Core module of the program under compilation is merged with the Core modules of all used libraries. The resulting big Core module is then translated to Grin.

4.2 Grin Transformations

In the Grin world, we take the opportunity to perform many optimizing transformations. Other transformations are designed to move from graph manipulation concepts (complete nodes that can be 'fetched', 'evaluated' and pattern matched for) to a lower level where single word values are moved and inspected in the imperative target language.

We first list all transformations in the order they are performed, and then discuss some issues that are tackled with the combined effort of multiple transformations.

- DropUnreachableBindings Drops all functions not reachable from main, either through direct calls, or through nodes that store a deferred or partially applied function. The transformation performs a provisional numbering of all functions, and creates a graph of dependencies. A standard graph reachability algorithm determines which functions are reachable from main; the others are dropped. This transformation is done as very first, because is drastically reduces program size: all unused functions from included libraries are removed.
- 2. *MergeInstance* Introduces an explicit dictionary for each instance declaration, by merging the default definitions of functions taken from class declarations. This is possible because we have the whole program available now (see discussion below).
- MemberSelect Looks for the selection of a function from a dictionary and its subsequent application to parameters. Replaces that by a direct call.
- 4. *DropUnreachableBindings* (again) Drops the now obsolete implicit constructions of dictionaries.
- Cleanup Replaces some node tags by equivalent ones: PApp 0, a partial application needing 0 more parameters, is changed into Fun, a simple deferred function; deferred applications of

- constructor functions are changed to immediate application of the constructor function.
- SimpleNullary Optimises nullary functions that immediately return a value or call another function by inlining them in nodes that encode their deferred application.
- 7. ConstInt Replaces deferred applications of integer2int to constant integers by a constant int. This situation occurs for every numeric literal in an Int context in the source program, because of the way literals are overloaded in Haskell.
- 8. *BuildAppBindings* Introduces bindings for *apply* functions with as many parameters as are needed in the program.
- GlobalConstants Introduces global variables for each constant found in the program, instead of allocating the constants locally.
- 10. Inline Inlines functions that are used only once at their call site.
- 11. *SingleCase* Replaces case expressions that have a single alternative by the body of that alternative.
- 12. EvalStored Do not do Eval on pointers that bind the result of a previous Store. Instead, do a Call if the stored node is a deferred call (with a Fun tag), or do a Unit of the stored node for other nodes.
- 13. ApplyUnited Do not perform Apply on variables that bind the result of a previous Unit of a node with a PApp tag. Instead, do a Call of the function if it is now saturated, or build a new PApp node if it is undersaturated.
- 14. SpecConst Specialize functions that are called with a constant argument. The transformation is useful for creating a specialized 'increment' function instead of plus 1, but its main merit lies in making specialized versions of overloaded functions, that is functions that take a dictionary argument. If the dictionary is a constant, specialization exposes new opportunities for the MemberSelect transformation, which is why SpecConst is iterated in conjunction with EvalStored, ApplyUnited and MemberSelect.
- 15. *DropUnreachableBindings* Drops unspecialized functions that may have become obsolete.
- NumberIdents Attaches an unique number to each variable and function name.
- 17. *HeapPointsTo* Does a 'heap points to analysis' (HPT), which is an abstract interpretation of the program in order to determine the possible tags of the nodes that each variable can refer to.
- 18. *InlineEA* Replaces all occurrences of *Eval* and *App* to equivalent constructs. Each *Eval* x is replaced by *FetchNode* x, followed by a *Case* distinction on all possible tag values of the node referred to by x, which was revealed by the HPT analysis. If the number of cases is prohibitively large, we resort to a *Call* to a generic *evaluate* function, that is generated for the purpose and that distinguishes all possible node tags. Each $App \ f \ x$ construct, that is used to apply an unknown function f to argument x, is replaced by a *Case* distinction on all possible PApp tag values of the node referred to by f.
- 19. *ImpossibleCase* Removes alternatives from *Case* constructs that, according to the HPT analysis, can never occur.
- LateInline Inlines functions that are used only once at their call site. New opportunities for this transformation are present because the *InlineEA* transformation introduces new *Call* constructs.
- 21. *SingleCase* (again) Replaces case expressions that have a single alternative by the body of that alternative. New opportunities for this transformation are present because the *InlineEA* transformation introduces new *Case* constructs.

- 22. DropUnusedExpr Removes bindings to variables if the variable is never used, but only when the expression has no side effect. Therefore, an analysis is done to determine which expressions may have side effects. Update and FFI expressions are assumed to have side effects, and Case and Seq expressions if one of their children has them. The tricky one is Call, which has a side effect if its body does. This is circular definition of 'has a side effect' if the function is recursive. Thus we take a 2-pass approach: a 'coarse' approximation that assumes that every Call has a side effect, and a 'fine' approximation that takes into account the coarse approximation for the body. Variables that are never used but which are retained because of the possible side effects of their bodies are replaced by wildcards.
- 23. *MergeCase* Merges two adjacent *Case* constructs into a single one in some situations.
- 24. LowerGrin Translates to a lower level version of Grin, in which variables never represent a node. Instead, variables are introduced for the separate fields, of which the number became known through HPT analysis. Also, after this transformation Case constructs scrutinize on tags rather than full nodes.
- 25. CopyPropagation Shortcuts repeated copying of variables.
- 26. SplitFetch Translates to an even lower level version of Grin, in which the node referred to by a pointer is not fetched as a whole, but field by field. That is, the FetchNode expression is replaced by a series of FetchField expressions. The first of these fetches the tag, the others are specialized in the alternatives of the Case expression that always follows a FetchNode expression, such that no more fields are fetched than required by the tag of each alternative.
- 27. *DropUnusedExpr* (again) Removes variable bindings introduced by *LowerGrin* if they happen not to be used.
- CopyPropagation Again shortcuts repeated copying of variables.

Simplification The Grin language has constructs for manipulating heap nodes, including ones that encode deferred function calls, that are explicitly triggered by an *Eval* expression. As part of the simplification, this high level construct should be decomposed in smaller steps. Two strategies can be used:

- tagged: nodes are tagged by small numbers, evaluation is performed by calling a special evaluate function that scrutinizes the tag, and for each possible Fun tag calls the corresponding function and updates the thunk;
- tagless: nodes are tagged by pointers to code that does the call
 and update operations, thus evaluation is tantamount to just
 jumping to the code pointed to by the tag.

The tagged approach has overhead in calling <code>evaluate</code>, but the tagless approach has the disadvantage that the indirect jump involved may stall the lookahead buffer of pipelined processors. Boquist proposed to inline the <code>evaluate</code> function at every occurrence of <code>Eval</code>, where for every instance the <code>Case</code> expression involved only contains those cases which can actually occur. It is this approach that we take in UHC.

This way, they high level concept of *Eval* is replaced by lower level concepts of *FetchNode*, *Case*, *Call* and *Update*. In turn, each *FetchNode* expression is replaced by a series of *FetchField* expressions in a later transformation, and the *Case* that scrutinizes a node is replaced by one that scrutinizes the tag only.

Abstract interpretation The desire to inline a specialized version of evaluate at every Eval instance brings the need for an anal-

ysis that, for each pointer variable, determines the possible tags of the node. An abstract interpretation of the program, known as 'heap points to (HPT) analysis' tries to approximate this knowledge. As preparation, the program is scanned to collect constraints on variables. Some constraints immediately provide the information needed (e.g., the variable that binds the result of a *Store* expression is obviously a pointer to a node with the tag of the node that was stored), but other constraints are indirect (e.g., the variable that binds the result of a *Call* expression will have the same value as the called function returns). The analysis is essentially a whole-program analysis, as actual parameters of functions impose constraints on the parameters.

The constraint set is solved in a fixpoint iteration, which processes the indirect constraints based on information gathered thus far. In order to have fast access to the mapping that records the abstract value for each variable, we uniquely number all variables, and use mutable arrays to store the mapping.

The processing of the constraint that expresses that x binds the result of $Eval\ p$ deserves special attention. If p is already known to point to nodes with a Con tag (i.e., values) then this is also a possible value for x. If p is known to point to nodes with a $Fun\ f$ tag (i.e., deferred functions), then the possible results for f are also possible values for f. And if f is known to point to nodes with an f apply tag (i.e., generic applications of unknown functions by f apply, then the possible results for f are also possible values for f. For a more detailed description of the algorithm, we refer to another paper [14].

HPT performance The HPT analysis must at least find all possible tags for each pointer, but it is sound if it reports a superset of these. The design of the HPT analysis is a tradeoff between time (the number of iterations it takes to find the fixed point) and accuracy. A trivial solution is to report (in 1 step) that every pointer may point to every tag; a perfect solution would solve the halting problem and thus would take infinite time in some situations.

We found that the number of iterations our implementation takes is dependent of two factors: the depth of the call graph (usually bounded by a dozen or so in practice), and the length of static data structures in the program. The latter surprised us, but is understandable if one considers the program

```
main = putStrLn (show (last [id, id, id, id, succ] 1))
```

where it takes 5 iterations to find out that 1 is a possible parameter of succ.

As for accuracy, our HPT algorithm works well for first-order functions. In the presence of many higher-order functions, the results suffer from 'pollution': the use of a higher-order function in one context also influences its result in another context. We counter this undesired behavior in several ways:

- instead of using a generic *apply* function, the *BuildAppBindings* transformation makes a fresh copy for each use by an *App* tag. This prevents mutual pollution of *apply* results, and also increases the probability that the *apply* function can be inlined later:
- we specialize overloaded functions for every dictionary that it is used with, to avoid the *App* needed on the unknown function taken from the dictionary;
- we fall back on explicitly calling *evaluate* (instead of inlining it) in situations where the number of possible tags is unreasonable large.

Instance declarations The basic idea of implementing instances is simple: an instance is a tuple (known as a 'dictionary') containing

all member functions, which is passed as an additional parameter to overloaded functions. Things are complicated, however, by the presence of default implementations in classes: the dictionary for an instance declaration is a merge of the default implementations and the implementations in the instance declaration. Worse, the class declaration may reside in another module than the instance declaration, and still be mutually dependent with it. Think of the Eq class, having mutually circular definitions of eq and ne, leaving it to the instance declaration to implement either one of them (or both).

A clever scheme was designed by Faxén to generate the dictionary from a generator function that is parameterized by the dictionary containing the default implementations, while the default dictionary is generated from a generator function parameterized by the instance dictionary [13]. Lazy evaluation and black holes make this all work, and we employ this scheme in UHC too. It would be a waste, however, now that we are in a whole-program analysis situation, not to try to do as much work as possible at compile time.

Firstly, we have to merge the default and instance dictionaries. In the Grin world, we have to deal with what the Core2Grin transformation makes of the Faxén scheme. That is:

- A 1-ary generator function *gfd* that, given a default dictionary, will generate the dictionary;
- A 0-ary function fd that binds a variable to a black hole, calls gfd, and returns the result
- A global variable d which is bound to a node with tag Fun fd.

We want to change this in a situation where d is bound directly to the dictionary node. This involves reverse engineering the definitions of d, fd and gfd to find the actual member function names buried deep in the definition of gfd. Although possible, this is very fragile as it depends on the details of the Core2Grin translation. Instead, we take a different approach: the definition of fd is annotated with the names of the member functions at the time when they are still explicitly available, that is during the EH2Core translation. Similarly, class definitions are annotated with the names of the default functions. Now the Grin.MergeInstance transformation can easily collect the required dictionary fields, provided that the Core.LiftDictFields transformation ensures they are available as top-level functions. The fd and gfd functions are obsolete afterwards, and can be discarded by a later reachability analysis.

Secondly, we hunt the program for dictionaries d (as constructed above) and selection functions s_k (easily recognizable as a function that pattern-matches its parameter to a dictionary structure and returns its kth field x_k). In such situations $Call \ s_k \ d$ can be replaced by $Eval \ x_k$. A deferred member selection, involving a node with tag $Fun \ s_k$ and field d, is dealt with similarly: both are done by the MemberSelect transformation.

Thirdly, as x_k is a dictionary field, it is a known node n. If n has a $Fun\ f$ tag, then $Eval\ x_k$ can be replaced by $Call\ f$, and otherwise it can be replaced by $Unit\ n$. This is done by the EvalStored transformation. The new $Unit\ that$ is exposed by this transformation can be combined with the App expression that idiomatically follows the member selection, which is what ApplyUnited does.

All of this only works when members are selected from a constant dictionary. Overloaded functions however operate on dictionaries that are passed as parameter, and member selection from a variable dictionary is not caught by *MemberSelect*. The constant dictionary appears where the overloaded function is called, and can be brought to the position where it is needed by specializing functions when they are called with constant arguments. This is done in the *SpecConst* transformation. That transformation is not only useful in the chain of transformations that together remove the dictionaries,

but also for the removal of other constant arguments, giving e.g. a 1-ary successor function as a specialization of plus 1. (If constant specialization is also done for string constants, we get many specializations of putStrLn).

The whole pack of transformations is applied repeatedly, as applying them exposes new opportunities for sub-dictionaries. Four iterations suffice to deal with the common cases (involving Eq, Ord, Integral, Read etc.) from the prelude.

The only situation where dictionaries cannot be eliminated completely, is where an infinite family of dictionaries is necessary, such as arises from the $Eq\ a\Rightarrow Eq\ [a]$ instance declaration in the prelude. We then automatically fall back to the Faxén scheme.

4.3 Silly Transformations

- 1. *InlineExpr* Avoids copying variables to other variables, if in all uses the original one could be used just as well (i.e., it is not modified in between).
- 2. *ElimUnused* Eliminates assignments to variables that are never used
- 3. *EmbedVars* Silly has a notion of function arguments and local variables. After this transformation, these kind of variables are not used anymore, but replaced by explicit stack offsets. So, this transformation does the mapping of variables to stack positions (and, if available, registers). In a tail call, the parameters of the function that is called overwrites the parameters and local variables of the function that does the call. The assignments are scheduled in such a way that no values are overridden that are still needed in assignments to follow.
- 4. GroupAllocs This transformation combines separate, adjacent calls to malloc into one, enabling to do heap overflow check only once for all the memory that is allocated in a particular function.

5. Conclusion

5.1 Code size

UHC is the standard materialization of a more general code base (the UHC framework, formerly known as EHC), from which increasingly powerful 'variants' of the compiler can be drawn, where independent experimental 'aspects' can be switched on or off. The whole source code base consists of a fairly exact 100.000 lines of code. Just over half of it is Attribute Grammar code, which of course has lots of embedded Haskell code in it. One third of the code base is plain Haskell (mostly for utility functions, the compiler driver, and the type inferencer), and one sixth is C (for the runtime system and a garbage collector).

In Figure 2 the breakdown of code size over various subsystems in the pipeline is shown. All numbers are in kilo-lines-of-code, but because of the total of 100.000 lines they can also be interpreted as percentages. Column 'UHC only' shows the size of the code that is selected by Shuffle for the standard compiler, i.e. the most powerful variant without experimental aspects. On average, 60% of the total code base is used in UHC. The rest is either code for low variants which is overwritten in higher variants, code for experimental aspects that are switched off in UHC, chunk header overhead, or comments that were placed outside chunks.

The fraction of code used for UHC is relatively low in the type inferencer (as there are many experimental aspects here), in the experimental backends like Java, Cil and LLVM (as most of them are switched off), and in the garbage collector (as it is not yet used: UHC by default uses the Boehm garbage collector [5, 4]).

subsystem	All variants and aspects				UHC only	
	AG	HS	C	total	total	fract.
utility/general	1.7	18.3		20.0	14.0	70%
Haskell	6.7	3.3		9.9	6.9	70%
EH	11.2	0.6		11.8	6.7	57%
EH typing	8.0	7.5		15.5	7.0	45%
Core	7.1	1.0		8.0	4.7	58%
ByteCode	2.1			2.1	1.7	82%
Grin	11.3	1.6		12.9	8.5	66%
Silly	2.8			2.8	2.6	93%
exp.backends	2.5	0.4		2.9	0.8	26%
runtime system			8.1	8.1	6.2	77%
garb.collector			6.0	6.0	0.7	11%
total	53.4	32.5	14.1	100.0	59.8	60%

Figure 2. Code size (in 1000 lines of code) of source files containing Attribute Grammar code (AG), Haskell code (HS) and C code (C), for various subsystems. Column 'all variants' is the total code base for all variants and aspects, column 'UHC' is the selection of the standard compiler, where 'fract.' shows the fraction of the full code base that is selected for UHC.

5.2 Methodological observations

Aspect-oriented organization UHC and its framework use an aspect-wise organization in which as much as possible is described by higher level domain specific languages from which we generate lower level implementations. UHC as a framework offers a set of compilers, thus allowing picking and choosing a starting point for play and experimentation. This makes UHC a good starting point for research, but debugging is also facilitated by it. A problem can more easily be pinpointed to originate in a particular step of the whole sequence of language increments; the framework then allows to debug the compiler in this limited context, with less interaction by other features.

The stepwise organization, where language features are built on top of each other, offers a degree of isolation. Much better would be to completely independently describe language features. However, this is hard to accomplish because language features often interact and require redefinition of parts of their independent implementation when combined. To do this for arbitrary combinations would be more complicated then to do it for a sequence of increments.

Testing can also be kept relatively simple this way. As long as an increment in features does not remove previous features or only changes the generated test output, tests for a previous step can still be reused and extended with new tests. In UHC this only fails when the presence of a Prelude is assumed; the testing framework is aware of this

The aspect-wise organization impacts all source code: AG code, Haskell code, C code, the build system, etc.. Implementing aspects as part of the used languages would be a major undertaking, as all languages then should be aware of aspects, and in a similar way. In UHC we have chosen to factor out aspect management and deal with it by preprocessing.

UHC as an experimentation platform — An obvious tension exists between UHC as a "full Haskell compiler" and a "nimble compiler for experimentation". Many seemingly innocent paragraphs of the Haskell language report have major impact on the implementation, making the implementation disproportional complex. Although this cannot be avoided, it can be isolated to a certain degree, which is what we hope to have accomplished using an aspect-wise approach. Although the chosen layering of language features and implementation techniques restricts the extent one can deviate from it for

experimentation, one can always select a minimal starting point in the sequence of compilers and build on top of that. When we add new functionality, we usually start by making it work in an early variant, and then gradually make it work for subsequent variants.

AG Design Patterns We tend to use various AG idioms frequently. For example, information is often gathered over a tree via a synthesized attribute, and subsequently passed back as an inherited attribute. This leads to a "cyclic program" when lazy code is generated from the AG description, or a 2-pass tree traversal when strict code is generated (after checking for absence of cycles).

Some idiomatic use is directly supported by the AG system. For example, transformations are expressed as attribute grammars with a single, specially designated, attribute declaration for a copy of the tree being walked over. The only thing that remains to be specified is where the transformed tree differs from the original.

The AG notation allows us to avoid writing much boilerplate code, similar to other tree traversal approaches [37, 36, 24]. The use of attributes sometimes also resembles reader, writer, and state monads. In practice, the real strength of the AG system lies in combining separately defined tree traversals into one. For example, the EH type analysis repeatedly builds environments for kinds, types, datatypes, etc. Combined with the above idiomatic use this easily leads to many passes over the EH tree; something we'd rather not write by hand using monads (and monad transformers) or other mechanisms more suitable for single-pass tree traversals!

However, not all idiomatic use is supported by AG. For example, the need to pattern match on subtrees arises when case analysis on abstract syntax trees must be done. Currently this must be programmed by hand, and we would like to have automated support for it (as in Stratego [37, 36]).

The use of intermediate languages UHC uses various intermediate languages and transformations on them. The benefit of this approach is that various compiling tasks can be done where it best fits an intermediate language and can be expressed as small, easy to understand, transformations independently from other tasks. Drawbacks are that some tasks have more than one appropriate place in the pipeline and sometimes require information thrown away in earlier stages (e.g. absence of types in Core).

The use of domain specific languages (DSL) We use various special purpose languages for subproblems: AG for tree traversals, Shuffle for incremental, aspect-wise, and better explainable development, Ruler for type systems. Although this means a steeper learning curve for those new to the implementation, in practice the DSLs we used and their supporting tools effectively solve an identifiable design problem.

5.3 Related work

Clearly other Haskell compilers exist, most notably GHC [26], which is hard if not impossible to match in its reliability and feature richness: UHC itself uses GHC as its main development tool.

Recently, JHC [27] and LHC [18] (derived from JHC) also take the whole-program analysis approach proposed by Boquist [8, 7] as their starting point. LHC in its most recent incarnation is available as a backend to GHC, and thus is not a standalone Haskell compiler.

Already longer available alongside GHC are Hugs [21] which was influential on Haskell as a language, NHC98 [38], and YHC [33] derived from NHC98, all mature Haskell 98 compilers with extensions. Helium [17] (also from Utrecht) does not implement full Haskell 98 but focuses on good error reporting, thereby being suitable for learning Haskell. We also mention HBC [1] (not maintained anymore) for completeness.

The distinguishing feature of UHC is its internal organization. UHC, in particular its internal aspect-wise organized framework, is designed to be (relatively) easy to use as a platform for research and education. In Utrecht students regularly use the UHC framework to experiment with. The use of AG and other tools also make UHC different from other Haskell compilers, most of them written in Haskell or lower level languages.

5.4 Future work

We have recently made a first public release of UHC [11]. In the near future we intend to add support for better installation, in particular the use of Cabal, and to add missing language features and libraries. On a longer time scale we will continue working on whole-program analysis, the optimizations allowed by it, add classical analyses (e.g. strictness), and improve the runtime system (switching to our own garbage collector). As we recently included the standard libraries, we will be able to run benchmark suites to compare the performance (code size, compilation time, run time) of each operation mode (bytecode interpreter, whole-program analysis) with each other and with other compilers. We welcome those who want to contribute in these or other areas of interest.

References

- L. Augustsson. The HBC compiler. http://www.cs.chalmers.se/~augustss/hbc/hbc.html, 1998.
- [2] R. Bird and O. de Moor. The algebra of programming. Prentice Hall, 1996.
- [3] R. S. Bird. Using Circular Programs to Eliminate Multiple Traversals of Data. Acta Informatica, 21:239–250, 1984.
- [4] H. Boehm. A garbage collector for C and C++. http://www.hpl.hp.com/personal/Hans_Boehm/gc/, 2006.
- [5] H. Boehm and M. Weiser. Garbage Collection in an Uncooperative Environment. Software Practice and Experience, pages 807–820, Sep 1988.
- [6] M. Bolingbroke and S. Peyton Jones. Types are calling conventions (submitted to Haskell Symposium 2009). 2009.
- [7] U. Boquist. Code Optimisation Techniques for Lazy Functional Languages, PhD Thesis. Chalmers University of Technology, 1999.
- [8] U. Boquist and T. Johnsson. The GRIN Project: A Highly Optimising Back End For Lazy Functional Languages. In Selected papers from the 8th International Workshop on Implementation of Functional Languages, 1996.
- [9] A. Dijkstra. Stepping through Haskell. PhD thesis, Utrecht University, Department of Information and Computing Sciences, 2005.
- [10] A. Dijkstra, J. Fokker, and S. D. Swierstra. The Structure of the Essential Haskell Compiler, or Coping with Compiler Complexity. In Implementation of Functional Languages, 2007.
- [11] A. Dijkstra, J. Fokker, and S. D. Swierstra. UHC Utrecht Haskell Compiler.http://www.cs.uu.nl/wiki/UHC, 2009.
- [12] A. Dijkstra and S. D. Swierstra.Ruler: Programming Type Rules. In Functional and Logic Programming: 8th International Symposium, FLOPS 2006, Fuji-Susono, Japan, April 24-26, 2006, number 3945 in LNCS, pages 30–46. Springer-Verlag, 2006.
- [13] K.-F. Faxén. A Static Semantics for Haskell. *Journal of Functional Programming*, 12(4):295, 2002.
- [14] J. Fokker and S. D. Swierstra. Abstract interpretation of functional programs using an attribute grammar system. In A. Johnstone and J. Vinju, editors, *Language Descriptions*, *Tools and Applications* (LDTA08), 2008.
- [15] GHC Team. The New GHC/Hugs Runtime System. http://citeseer.ist.psu.edu/marlow98new.html, 1998.

- [16] Haskell' Committee. Haskell Prime. http://hackage.haskell.org/trac/haskell-prime/, 2009.
- [17] B. Heeren, A. v. IJzendoorn, and J. Hage.Helium, for learning Haskell. http://www.cs.uu.nl/helium/, 2005.
- [18] D. Himmelstrup, S. Bronson, and A. Seipp. LHC Haskell Compiler. http://lhc.seize.it/, 2009.
- [19] ISO. Common language infrastructure (ISO/EIC standard 23271). ECMA. 2006.
- [20] M. P. Jones. Typing Haskell in Haskell. In Haskell Workshop, 1999.
- [21] M. P. Jones. Hugs 98. http://www.haskell.org/hugs/, 2003.
- [22] D. Knuth. Semantics of context-free languages. Mathematical Systems Theory, 2(2):127–145, 1968.
- [23] D. Knuth. Literate Programming. Journal of the ACM, (42):97–111, 1984.
- [24] R. Lämmel and S. Peyton Jones. Scrap your boilerplate: a practical design pattern for generic programming. In *Types In Languages Design And Implementation*, pages 26–37, 2003.
- [25] C. Lattner and V. Adve. LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation. In *Proceedings of the* 2004 International Symposium on Code Generation and Optimization (CGO'04), Palo Alto, California, Mar 2004.
- [26] S. Marlow and S. Peyton Jones. The Glasgow Haskell Compiler. http://www.haskell.org/ghc/, 2004.
- [27] J. Meacham. Jhc Haskell Compiler. http://repetae.net/computer/jhc/, 2009.
- [28] S. Peyton Jones. Compiling Haskell by program transformation: a report from the trenches. In *European Symposium On Programming*, pages 18–44, 1996.
- [29] S. Peyton Jones. Tackling the Awkward Squad: monadic input/output, concurrency, exceptions, and foreign-language calls in Haskell. In Engineering theories of software construction, Marktoberdorf Summer School, 2002.
- [30] S. Peyton Jones. Haskell 98, Language and Libraries, The Revised Report. Cambridge Univ. Press, 2003.
- [31] S. Peyton Jones and S. Marlow. Secrets of the Glasgow Haskell Compiler inliner. *Journal of Functional Programming*, pages 393– 434, 2002.
- [32] S. Peyton Jones and E. Meijer. Henk: A Typed Intermediate Language. In Workshop on Types in Compilation, 1997.
- [33] T. Shackell, N. Mitchell, A. Wilkinson, et al. YHC York Haskell Compiler.http://haskell.org/haskellwiki/Yhc, 2009.
- [34] S. D. Swierstra, P. Azero Alocer, and J. Saraiva. Designing and Implementing Combinator Languages. In 3rd Advanced Functional Programming, number 1608 in LNCS, pages 150–206. Springer-Verlag, 1999.
- [35] M. Viera, S. D. Swierstra, and W. S. Swierstra. Attribute grammars fly first class: How to do aspect oriented programming in haskell. In *International Conference on Functional programming (ICFP '09)*, New York, NY, USA, 2009. ACM Press.
- [36] E. Visser. Stratego: A language for program transformation based on rewriting strategies. System description of Stratego 0.5. In A. Middeldorp, editor, *Rewriting Techniques and Applications* (RTA'01), number 2051 in LNCS, pages 357–361. Springer-Verlag, 2001.
- [37] E. Visser. Stratego Home Page. http://www.program-transformation.org/Stratego/WebHome, 2005.
- [38] York Functional Programming Group. NHC98 Haskell Compiler. http://haskell.org/nhc98/, 2007.