Abstract Syntax Graphs for

Bruno C. d. S. Oliveira

National University of Singapore oliveira@comp.nus.edu.sg

Abstract

cific languages (EDSLs) using abstract syntax graphs (ASGs). The purpose of this representation is to deal with the important problem of defining operations that require observing or preserving sharing and recursion in EDSLs in an expressive, yet easy-to-use way. In contrast to more conventional representations based on abstract syntax trees, ASGs represent sharing and recursion explicitly as binder constructs. We use a functional representation of ASGs based on structured graphs, where binders are encoded with parametric higher-order abstract syntax. We show how adapt to this representation to well-typed ASGs. This is especially useful for EDSLs, which often reuse the type system of the host language. We also show an alternative class-based encoding of (well-typed) ASGs that enables extensible and modular well-typed EDSLs while allowing the manipulation of sharing and recursion.

This paper presents a representation for embedded domain spe-

Categories and Subject Descriptors D.3.2 [Programming Languages]: Language Classifications—Functional Languages

General Terms Languages

Keywords Observable Sharing, DSLs, Graphs, Haskell.

1. Introduction

A domain-specific language (DSL) is a programming language targeted at a particular problem domain. DSLs offer a vocabulary, language constructs and a semantics crafted for that domain.

An embedded DSL (EDSL) [14] is a DSL that is implemented by reusing various elements of a (general-purpose) host language (such as the syntax, type-checker, or binding constructs). While being somewhat less flexible than writing a dedicated compiler or interpreter for a DSL, the embedded approach greatly reduces the cost of the implementation. Furthermore, integration with the host language comes for free, and mixing the DSL with the host language or other DSLs is easy.

Representation of the syntax of an EDSL in the host language are typically positioned between two extremes: a *shallow* embedding provides a very thin layer over the host language, implementing the DSL constructs directly by their semantics. As a result, there is no support for inspecting the syntax and manipulations of the DSL programs are difficult. *Deep* embeddings solve this problem by making the syntax of the DSL explicit – usually as an *abstract*

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.

PEPM'13, January 21–22, 2013, Rome, Italy. Copyright © 2013 ACM 978-1-4503-1842-6/13/01...\$15.00

Domain Specific Languages

Andres Löh Well-Typed LLP andres@well-typed.com

given by defining interpretation functions over the AST, and transformations of the AST prior to interpretation are possible. The AST approach is well supported by functional languages such as Haskell and has been used to implement several EDSLs [4, 19].

syntax tree (AST). One or several semantics of the DSL can be

In many real-life EDSLs, preserving and observing sharing and recursion are essential for implementing domain-specific transformations and optimizations. However, ASTs need to be complemented with *explicit* environments to allow transformations that rely on observing sharing or recursion. Furthermore, doing so, we suddenly need to keep track of names and binding, forcing us to a lower level of programming, where we have to worry about problems such as avoiding name capture in substitutions or preventing dangling references.

Abstract Syntax Graphs This paper suggests using an abstract syntax graph (ASG) representation for EDSLs. ASGs make it easy to guarantee that terms are well-scoped. They allow the observation and preservation of sharing and recursion. Furthermore, functions on ASGs can be defined in a natural way, using pattern matching.

Cook's structured graphs [20]. Such structured graphs offer a generic purely functional representation of cyclic structures in pure functional languages such as Haskell. Structured graphs use binders, represented using Chlipala's parametric higher-order abstract syntax (PHOAS) [7], to model cycles and sharing.

Related Work With ASTs, a possible approach to help with the

issues regarding the management of explicit environment and generation of fresh labels [4] is to use *monads* [26]. However, while

Technically speaking ASGs are realized using Oliveira and

monads can make name management more bearable, they cannot completely hide the fact that we have to work on a low level, and monads alone cannot ensure the well-scopedness of a program.

Sometimes, one would like to ensure not only well-scoped, but also well-typed expressions in an EDSL. Guaranteeing well-typedness becomes even harder in the presence of explicit environments. Both the ASTs and the respective environments need to be enriched with additional type and binding information. Examples of such approaches include well-typed and well-scoped analysis and transformation of grammars, which have been a hot topic recently [2, 3, 9]. Baars et al. [2, 3] use well-typed ASTs in combination typed references and typed environments in their typed transformations. The relationship between a reference and an environment is statically enforced in a similar way to well-scoped/typed

type-level machinery and several Haskell extensions.

Another option is to represent sharing and recursion implicitly, by relying on the sharing of the host language. This is great from the usability point of view because we can reuse the host language syntax to create sharing. Unfortunately, this is usually too fragile or precludes the possibility of observation. To overcome the need for observing sharing and recursion in implicit representations, it

is possible to use pointer or reference equality. This approach has

de Bruijn indices [1]. All this infrastructure relies on sophisticated

defining operations by working with references is fragile and prone to errors. In contrast, our ASGs are completely functional and avoid the need for observing and comparing pointers.

More recently, both Devriese et al. [9, 10] and Kiselyov [16] have proposed the use of recursive binders to implement explicit sharing using a type-class based representation. However a problem of such type-class based representations is that pattern matching is not supported. All operations are essentially defined as folds, making it difficult to define many transformations and optimizations that rely on both observing sharing and use more complex recursion patterns that would be most naturally expressed using nested pattern matching.

*Contributions** The contributions of this paper 1 are:

been used in many DSL implementations [8, 11, 18]. However, the use of references breaks *referential transparency* and significantly complicates reasoning [25]. In a language such as Haskell, we would then be forced to use monadic interfaces. Furthermore

lieve that the ASG representation is particularly valuable to EDSL developers, providing the best solution to date to the problem of observing sharing and recursion.

We make a case for the use of these techniques and we also fill in some gaps that are not covered in the earlier works on PHOAS

ASGs for EDSLs Neither Oliveira and Cook nor Chlipala considered the application of PHOAS and ASGs to EDSLs. Yet, we be-

and ASGs, that are of relevance in the context of EDSLs.

Well-typed ASGs We show how to represent deeply embedded well-typed terms with ASGs. Well-typed terms allow DSL designers to reuse (parts of) the type system of the host language for the type system of the DSL.

Oliveira and Cook develop techniques only for representing untyped abstract syntax using structured graphs. Chlipala's original work on PHOAS does cover a form of well-typed abstract

sive binders and their mutually recursive generalization. While well-typed encodings of simple recursive binders are relatively straightforward to encode, encoding well-typed mutually recursive binders elegantly requires more work. To solve this problem we use *typed lists*: a generalization of both heterogeneous and homogeneous lists. With typed lists not only can we deal with well-typed mutually recursive binders, but also enforce certain size invariants that the untyped representation does not statically enforce

syntax using dependent types in the Coq theorem prover. However, conventional functional languages like Haskell do not have full-blown dependent types (although recent extensions get us very close to that [28]), so these techniques have to be adapted to use GADTs [23] instead. Also, Chlipala does not cover recur-

mutually recursive binders, but also enforce certain size invariants that the untyped representation does not statically enforce.

Extensible and modular ASGs The issue of modularity is not considered neither by Oliveira and Cook nor by Chlipala. A popular representation of EDSLs that deals with this problem uses a class-based (typed) Church encoding representation [5, 12, 13, 22] (although this representation makes the definition of operations that use nested pattern matching more difficult). We show that ASG-based techniques can be adapted to a class-based representation, similar to the representations proposed by Devriese et al. [9, 10] and Kiselvoy [16]. However, differently from these approaches, we

2. Sharing in EDSLs

In this section, we use a small example language to demonstrate

encoding of mutually recursive binders using typed lists.

the differences between shallow and deep embeddings as well as the issues with representing sharing.

use a PHOAS-based representation of binders and provide a simple

 $^{^1{\}rm The}$ code for this paper is available at http://ropas.snu.ac.kr/~bruno/papers/ASGs.zip.

In order to keep the examples as small as possible, we use an EDSL with just two constructs: the constant one, and a binary addition operator. The Haskell interface of our DSL is:

one :: Expr (\oplus) :: Expr \to Expr \to Expr

The primary semantics we are interested in is evaluation:

data Expr -- abstract

eval :: Expr ightarrow Int

We are now going to contrast a shallow embedding with a deep embedding for this language.

Shallow embedding A shallow embedding for this language is:

type Expr = Int one = 1

 $(\oplus) = (+)$ eval = id

We use the type Int as the representation of the expression type. Building a term in the expression language evaluates it automati-

cally. The evaluation function eval is then just the identity function.

The shallow approach is appealing because it is so simple. Constructing terms in the DSL is as easy as constructing Haskell terms.

We even inherit many features from the host-language Haskell. For example, we can use a Haskell function to generate a term in our DSL, as shown on the left hand side of Figure 1.

The term tree₁ n (Figure 1) describes a binary tree of additions, with occurrences of one in the leaves. The function tree₁ is recur-

pression DSL, yet they are available to us via the embedding into Haskell.

The use of sharing is essential here for efficient evaluation of the term. Without sharing, tree₁ n would contain exponentially many

sive, and it makes use of sharing via let. Both recursion and sharing are properties we do not have available in the interface of our ex-

additions and constants in n. By using sharing, the term is internally represented as a graph of just linear size. The identifier shared is bound to an Expr represented as an Int, and even though shared is

being used twice, it is being evaluated only once. The evaluation

of eval (tree₁ 2) is sketched on the left hand side of Figure 2. Note how 1 + 1 is evaluated only once, and its result (2) is shared.

However, shallow embeddings come at a price. We are committing to a specific semantics – in this case, evaluation. Often, that is

too limited in practice. We may want to do other things with expressions: for example, show the original term via a function

text::Expr → String

or transform the expression into a different (perhaps optimized) form, or translate the expression into a different language with a different set of constructs available. With a shallow embedding, we are out of luck. Our implementation picks one semantics and once we construct a term, we interpret the expression according to that

we construct a term, we interpret the expression according to that semantics, losing the original structure of the expression.

Deep embedding A deep embedding solves this problem:

data Expr = One | Add Expr Expr

 (\oplus) = Add eval One = 1 eval (Add e₁ e₂) = eval e₁ + eval e₂

one = One

stract syntax. A value of type Expr corresponds to the abstract syntax tree of a term in our DSL. We thus retain the structure of the terms we construct and can interpret them in various ways. We can,

We now choose to represent the language constructs by their ab-

tree, 0 = one tree, 0 = one tree, 0 = one tree, n = let shared = tree (n − 1) in shared⊕ shared tree either using Haskell's implicit sharing (left) or explicit sharing in our DSL (right)

trees :: Int → Expr

eval (tree, 2)

= eval (let shared = tree; (2 - 1) in Add shared shared)

```
= let shared = let shared' = tree (1 - 1) in shared' + shared'
                                                                           = let shared = tree; (2 - 1) in eval shared + eval shared
  in shared + shared
                                                                           = let shared = let shared' = tree; (1 - 1) in Add shared' shared'
- let shared - let shared' - 1 in shared' + shared'
                                                                              in eval shared + eval shared
  in shared + shared
                                                                           = let shared' = tree, (1-1)
= let shared = 1 + 1 in shared + shared
                                                                              in (eval shared' + eval shared') + (eval shared' + eval shared')
= let shared = 2 in shared + shared
                                                                           = let shared' = One
                                                                              in (eval shared' + eval shared') + (eval shared' + eval shared')
= 4
                                                                           =(1+1)+(1+1)
                                                                           = 2 + 2
                                                                           -4
```

Figure 2. Contrasting evaluation of eval (tree, 2) using both the shallow (left) and deep (right) embedding

for example, evaluate it as shown in the definition of eval above,

but we can also show it in textual form:

free :: Int → Expr

aval (tree, 2)

= let shared = tree; (2 - 1) in shared + shared

text :: Expr \rightarrow String

```
text One = "1"

text (Add e_1 e_2) = "("++ text e_1 ++ " + "+ text e_2 ++ ")"
```

In a similar way, we could define additional interpretation functions such as an optimizer or a translator to a different language. Typically, the interpretation functions are *folds* (also known as *catamorphisms*), i.e., functions that traverse the structure of the underlying input datatype (here Expr) closely and recurse exactly where we encounter a recursive subterm in the datatype definition.

However, the greater flexibility comes at a price. Consider tree₁ again, defined exactly as before (that is, the tree₁ definition in the left side of Figure 1). The identifier shared now is a term of the datatype Expr, no longer of type Int. If we evaluate the term tree₁ n

the sharing. The term will take exponentially long to evaluate (or to show, or to transform). The evaluation of eval (tree, 2) in the deep setting is sketched on the right side of Figure 2. Note how the pattern matching in eval destroys the sharing introduced by let, and how 1 + 1 is evaluated twice. Haskell's let still allows us to construct implicitly shared terms

using eval, we traverse the structure of the Expr, thereby destroying

of type Expr, but this sharing is not observable and is also quite fragile. Traversing such an implicitly shared term using any interpretation function will destroy all sharing.

Explicit sharing A solution is to make sharing explicit in the embedded language. This will enable us to observe and preserve

the sharing that we wish to have in a term in a robust way. It is quite clear that we need to add a let-like construct, but there is quite some design flexibility in the detail. We would like to avoid having to deal with names, binding and substitution ourselves, as this is tedious and error-prone, and would make the DSL much more tricky to use or at least to implement. One promising approach to model binding in the embedded

language is higher-order abstract syntax (HOAS) [24]. With HOAS the function space of the implementation language Haskell is used in order to express a shared term in the embedded language:

data Expr = One | Add Expr Expr | Let Expr (Expr → Expr)

We no longer have to use Haskell's let in order to express sharing in the embedded language. Next to one and (\oplus) (that can be defined

89 as before) we have to augment the interface of our language with an explicit sharing construct:

```
let_{::} Expr \rightarrow (Expr \rightarrow Expr) \rightarrow Expr
let = Let
```

tree, called tree, is shown on the right side of Figure 1. However there is a problem: How do we extend the evaluator to cover the case for Let? Here is an attempt:

We have to adapt the construction of shared terms to use this explicit sharing construct. The resulting modification of function

eval (Let $e_1 e_2$) = **let** shared = eval e_1 **in** eval (e_2 (... shared))

We would like to feed the evaluated shared shared expression to e₂, but it has the wrong type! The body of the Let expects an Expr, but we have an Int. At the position of ..., we need a function that

can quote the interpreted term back into the original language [17]. Alternatively, we have to add another constructor to Expr, because the existing constructors are not really expressive enough (we have One, but not arbitrary integer literals). Note that other interpretation

functions such as text would need other quotation functions. But before we delve too deep into this issue, we should point out another problem with higher-order abstract syntax: the space of type Expr o Expr is too large. In order to express binding faith-

fully, we want the syntactic shape of the resulting expression to be independent of the expression being shared. However, a Haskell function of Expr → Expr allows us to plug in functions that caseanalyze the incoming value and return different expressions depending on the outcome of that analysis.

Abstract syntax graphs Making sharing explicit means that the

abstract syntax representation becomes a graph rather than a tree. Although our effort to use HOAS to model ASGs has some problems, Oliveira and Cook [20] have shown a functional representa-

tion of graphs that solves these problems. The idea is to use parametric higher-order abstract syntax (PHOAS) [7] instead of HOAS

to model binders.

data Expr a =	= One	Add ((Expr a)	(Expr	a)
	Var a	Let ((Expr a)	$(a \rightarrow$	Expr a)

With PHOAS the whole expression datatype is now parameterized by the type of shared expressions a. We have two new constructors compared to our original type, one for variables that embeds a value of type a in Expr, and one for Let. The body of the Let now receives a variable of type a rather than a value of type Expr.

If we now require expressions in our language to make no as-

sumption about the variables, i.e., to be *polymorphic* in a, then (unlike HOAS) we cannot analyze the shared expression. Furthermore, Var serves as a generic way to quote intermediate results of interpretation functions. We can thus make the following definition for closed expressions, i.e., expressions with no free variables:

type ClosedExpr =
$$\forall$$
a.Expr a

We use ClosedExpr to explicitly refer to closed terms in our DSL and Expr a to construct terms or write interpreter functions.

We define one and (\oplus) as before:

one = One
$$(\oplus)$$
 = Add

In addition, we define a function let_ that wraps Let:

let_::Expr a
$$\rightarrow$$
 (Expr a \rightarrow Expr a) \rightarrow Expr a
let_e₁ e₂ = Let e₁ ($\lambda x \rightarrow$ e₂ (Var x))

The PHOAS underpinning guarantees that we cannot do anything with the argument we obtain in the body of the Let but to use it as a variable. But having to invoke Var explicitly at every use site is somewhat tedious – the wrapper performs this work for us.

With these definitions in place, we can define our explicitly shared tree_E function again. It looks just like the definition on the right side of Figure 1, but its type becomes $Int \rightarrow ClosedExpr$. We now have the choice whether to use Haskell's host-language let construct while doing meta-programming by writing a term like on the left side of Figure 1, or if we explicitly want to express sharing in the embedded language using let like on the right side.

Preserving sharing The evaluator can now be defined as follows:

eval:: Expr Int \rightarrow Int

```
\begin{array}{ll} \text{eval One} &= 1 \\ \text{eval } (\text{Add } \mathbf{e}_1 \ \mathbf{e}_2) = \text{eval } \mathbf{e}_1 + \text{eval } \mathbf{e}_2 \\ \text{eval } (\text{Var n}) &= \mathbf{n} \\ \text{eval } (\text{Let } \mathbf{e}_1 \ \mathbf{e}_2) &= \text{eval } (\mathbf{e}_2 \ (\text{eval } \mathbf{e}_1)) \end{array}
```

text::ClosedExpr \rightarrow String

The interpreter expects an Expr Int – it thus assumes that variables are of type integer for the purpose of evaluating an expression. However, a ClosedExpr is polymorphic in the variable type, so it will naturally be accepted by eval. In the Var case, we find an integer and can return it. In the Let case, we have to provide an integer for the value of the bound variable: we pass eval e₁. Note that this achieves sharing, because lambda-bound terms in Haskell are automatically shared. Therefore calling eval (tree_E 30) now will return the result 1073741824 almost immediately.

tation functions for expressions. Here is a function that computes a textual representation of the given term. Here, rather than *preserving* the sharing, we are interested in *observing* it:

Observing sharing Furthermore, it is easy to write other interpre-

```
text e = go e 0

where

go:: Expr String \rightarrow Int \rightarrow String

go One

_= "1"

go (Add e<sub>1</sub> e<sub>2</sub>) c =

"("++go e<sub>1</sub> c++" + "++go e<sub>2</sub> c++")"

go (Var x)

_= x

go (Let e<sub>1</sub> e<sub>2</sub>) c =

"(let "++v++" = "++go e<sub>1</sub> (c+1)++

" in "++go (e<sub>2</sub> v) (c+1)++")"

where v = "v"++show c
```

Inlining As a final example, let us look at a transformation that removes explicit sharing again, effectively inlining all let-bound

text (tree_F 2) vields

variables:

inline One = One inline (Add $e_1 e_2$) = Add (inline e_1) (inline e_2) inline (Var x) = x inline (Let $e_1 e_2$) = inline (e_2 (inline e_1))

inline:: Expr (Expr a) \rightarrow Expr a

This operation produces the original expression, but unfolds Let constructs. For the purposes of inline, variables are themselves expressions. For text (inline (tree_E 2)), we obtain "((1 + 1) + (1 + 1))"

Summary We have shown that there are situations where we need

a let-construct rather than unfolding the expression. Evaluating

"(let v0 = (let v1 = 1 in (v1 + v1)) in (v0 + v0))"

again, and eval (inline (tree_E 30)) takes forever to compute.

to observe or preserve sharing in an embedded DSL. Preserving sharing may be needed for performance reasons (as in the tree example), or it may be needed for operations that inspect shared terms and treat them in a particular way (as in the text example). PHOAS offers a safe yet convenient way to make sharing explicit and encode ASGs. The user can reuse Haskell's own scoping rules and does not have to worry about managing names. Differ-

To this end, we extend our example language with a few new constructs. For now, let us move from the constant "one" to allowing arbitrary integer literals, add a construct for checking if a term

In this section, following Oliveira and Cook [20], we will extend the solution to sharing presented in Section 2 to recursive and

is equal to "zero", and add lambdas and application:

sis on bound variables is forbidden.

(Mutual) recursion

mutually recursive bindings.

ently from classic HOAS encoding, terms that perform case analy-

type ClosedExpr = $\forall a$.Expr a data Expr a = Lit Int | Add (Expr a) (Expr a) IfZero (Expr a) (Expr a) (Expr a) Var a | Let (Expr a) $(a \rightarrow Expr a)$ Lam $(a \rightarrow Expr a) \mid App (Expr a) (Expr a)$

exactly as before. In IfZero, we take a condition, a then-part and an else-part. Variables (Var) and Let are unchanged. A lambda (Lam) is a binding construct. It therefore takes a function of type $a \rightarrow Expr a$ in the same way as the body of Let. Application (App) takes a

The constructor Lit takes an arbitrary integer literal. Addition is

function and an argument. We define a few "smart constructors" to facilitate constructing terms again:

 (\oplus) = Add $(\odot) = App$ $let_e_1 e_2 = Let e_1 (\lambda x \rightarrow e_2 (Var x))$ $lam_e = Lam(\lambda x \rightarrow e(Var x))$

Evaluation Let us look at how to extend the evaluator. We no

evaluate to integers. Instead, terms of our language now have a type τ where the type language is as follows: $\tau ::= \text{Int} \mid \tau \to \tau$

We have some flexibility encoding the type system when we embed the language: we can encode the types of the terms dynamically, and allow the language to represent ill-typed terms that will fail at run-time; or we can use Haskell's type system to enforce that terms in the language must be well-typed. Both settings have some merit. We will therefore look at the dynamic approach here and deal with

longer have the luxury that all terms of our embedded language

eval (Add
$$e_1 e_2$$
) = add (eval e_1) (eval e_2)
eval (IfZero $e_1 e_2 e_3$) = ifZero (eval e_1) (eval e_2) (eval e_3)

Functions are represented as Haskell functions in this simple setting – we might move to a representation using an explicit closure using an environment in a larger setting. The evaluator changes slightly

the static encoding of the types later, in Section 5.

The result of evaluation is now a tagged value:

data Value = N Int | F (Value → Value)

an environment in a larger setting. The evaluator changes slightly as a consequence, and now looks as follows:

eval:: Expr Value \rightarrow Value

eval (Lit i) = N i

eval (Add e₁ e₂) = add (eval e₁) (eval e₂)

```
\begin{array}{ll} \text{eval (Var x)} &= x \\ \text{eval (Let } e_1 e_2) &= \text{eval (} e_2 \text{ (eval } e_1)\text{)} \\ \text{eval (Lam } e) &= F \left(\lambda v \rightarrow \text{eval (} e v\text{)}\right) \\ \text{eval (App } e_1 e_2) &= \text{app (eval } e_1) \text{ (eval } e_2) \end{array}
```

We now have to tag values whenever we produce them, such as in the cases for Lit and Lam. For constructors such as Add, IfZero and App we write wrapper functions that check (at run time) whether the arguments have the correct types and throw an error if not:

```
add (N m) (N n) = N (m + n)
```

app (Ff)v' = fvOf course, we could also define a monadic evaluator that would be a total function and raturn Maybe Value instead of Value

a total function and return Maybe Value instead of Value. Here is a small example:

 $\begin{array}{l} \text{let}_(\text{lam}_(\lambda x \to x \oplus \text{Lit} (-1))) \ (\lambda \text{dec} \to \\ \text{let}_(\text{lam}_(\lambda f \to \text{lam}_(\lambda x \to \\ f \odot (f \odot x)))) \ (\lambda \text{twice} \to \\ \text{(twice} \odot \text{twice} \odot \text{dec} \odot \text{Lit} \ 10))) \end{array}$

ifZero (N n) v_1 $v_2 = if$ n == 0 then v_1 else v_2

This expression encodes the term

let dec x = x - 1

example =

 $\begin{aligned} &\text{twice f } x = f\left(fx\right) \\ &\text{in twice twice dec } 10 \end{aligned}$ Note that the two uses of twice are at different types. Evaluating the

expression eval example yields N 6 as expected.

Recursion Recursion is simple to add, by introducing an additional constructor that represents fixed points:

data Expr $a = \dots$ -- as before \mid Mu $(a \to \text{Expr a})$ This binding construct is very similar to Lam. Both constructs

introduce a bound variable that scopes over the entire body of the expression.

The idea is that using Mu, we can encode a recursive function such as multiplication (in terms of addition) as follows:

 $\begin{aligned} &\text{mu_e} = \text{Mu } (\lambda x \to e \; (\text{Var } x)) \\ &\text{mul :: ClosedExpr} \\ &\text{mul} = \text{lam_} (\lambda m \to \text{mu_} (\lambda \text{rec} \to \text{lam_} (\lambda n \to \text{lfZero n } (\text{Lit } 0) \; (\text{m} \oplus (\text{rec} \odot (\text{n} \oplus \text{Lit } (-1))))))) \end{aligned}$

The evaluator must of course be adapted as well:

eval (Mu e) = fix ($\lambda v \rightarrow \text{eval (e } v)$)
The new case maps Mu to Haskell recursion using the fix function:

91

fix :: $(a \rightarrow a) \rightarrow a$ fix f = **let** r = f r **in** r

eval :: ClosedExpr \rightarrow Value eval ... = ... -- as before

As we did in Section 2, we can also write other semantic functions on our DSL such as a function text to display the expres-

Using the let here for the result introduces additional sharing.

sion. Semantic functions can now observe and preserve recursion as needed.

It is also possible to define a recursive let-construct in terms of

Let and Mu: letrec:: (Expr $a \to Expr a$) \to (Expr $a \to Expr a$) $\to Expr a$

letrec:: (Expr a \rightarrow Expr a) \rightarrow (Expr a \rightarrow Expr a) \rightarrow Expr a letrec e₁ e₂ = Let (Mu ($\lambda x \rightarrow e_1$ (Var x))) ($\lambda x \rightarrow e_2$ (Var x))

Mutually recursive definitions The Mu construct is sufficient for expressing simple recursion, but we cannot easily express the definition of several mutually recursive bindings. For languages with

an expressive internal structure we might be able to encode mutual recursion in terms of simple recursion within the DSL, but we want our techniques to be widely applicable and not impose strong requirements on the DSLs.

When defining mutually recursive definitions we need to bind several variables at once (one for each mutually recursive definition). As an example, consider the following Haskell term: let dec x = x - 1

even x t e = if x = 0 then t else odd (dec x) t eodd x t e = if x = 0 then e else even (dec x) t ein even 4 1 0

number is even, the first continuation is returned, if it is odd, then the second continuation is returned instead. The given call returns 1, because 4 is even.

The functions even and odd are mutually recursive, and both depend on dec. This kind of mutually recursive binding is com-

The function even takes a number and two continuations. If the

monplace in a language like Haskell.

To deal with mutually recursive bindings, we add a new constructor called LatBac:

structor called LetRec:

data Expr a = ... -- as before

$$\vdash \mathsf{LetRec}\;([\mathsf{a}] \to [\mathsf{Expr}\,\mathsf{a}])\;([\mathsf{a}] \to \mathsf{Expr}\,\mathsf{a})$$

refer to each of the others. So all declarations are parameterized by a list of inputs. The body also can refer to each of the bindings, therefore it is parameterized over the same list. The type system cannot express the intuition that all three lists that occur in the type

We have a list of declarations now. Each of the declarations can

above are supposed to have the same length. We will be able to make this precise in Section 5.

We also define a wrapper that applies Var to all the variables:

$$\begin{array}{c} \mathsf{letrec}_{::} ([\mathsf{Expr}\,\mathsf{a}] \to [\mathsf{Expr}\,\mathsf{a}]) \to \\ \qquad \qquad ([\mathsf{Expr}\,\mathsf{a}] \to \mathsf{Expr}\,\mathsf{a}) \to \mathsf{Expr}\,\mathsf{a} \\ \mathsf{letrec}_\,\mathsf{es}\,\mathsf{e} = \mathsf{LetRec}\,(\lambda \mathsf{xs} \to \mathsf{es}\,(\mathsf{map}\,\mathsf{Var}\,\mathsf{xs})) \\ \qquad \qquad (\lambda \mathsf{xs} \to \mathsf{e}\,\,(\mathsf{map}\,\mathsf{Var}\,\mathsf{xs})) \end{array}$$

[lam_
$$(\lambda x \to x \oplus \text{Lit} (-1))$$

, lam_ $(\lambda x \to \text{lam}_-(\lambda t \to \text{lam}_-(\lambda e \to \text{lfZero } x t \text{ (odd } \odot (\text{dec } \odot x) \odot t \odot e))))$
, lam_ $(\lambda x \to \text{lam}_-(\lambda t \to \text{lam}_-(\lambda e \to \text{lam}_-(\lambda e)))$

```
If Zero x e (even \odot (dec \odot x) \odot t \odot e))))
])
(\lambda[dec, even, odd] \rightarrow even \odot Lit 4 \odot Lit 1 \odot Lit 0)
```

The only slightly tricky point is that we need to delay the pattern match on the list of variables in the first argument to letrec_ (using \sim), because in an interpretation function, Haskell will not be able to determine the number of elements in this list before looking at the body of the lambda.

We can extend an interpretation function such as the evaluator to cope with the presence of LetRec as follows:

```
eval :: ClosedExpr \rightarrow Value
eval ... = ... -- as before
eval (LetRec es e) = eval (e (fix (map eval \circ es)))
```

Reusing native let syntax It can be argued that despite the advantages of using explicit sharing, it is still less convenient to use let_ or letrec_ than to use Haskell's native let construct.

Many EDSLs therefore use Haskell's **let**, but recover the sharing information by inspecting the internal representation of the term, using an impure function. The function reifyGraph, from the data-reify package [11], provides such functionality. This function returns a graph representing subterms using numbers – a representation that is neither particularly safe nor directly suitable for further computations.

We can combine reification with our ASG approach. We start with arithmetic expressions with just literals and addition:

```
data Expr_D = Lit_D Int | Add_D Expr_D Expr_D
```

The goal is to convert an implicitly shared term such as tree 3 (using tree 1 from Figure 1 with type Int \rightarrow Expr_D, with obvious definitions of one and \oplus) into an explicitly shared term of type Expr. In order to be able to use data-reify on terms of type Expr_D, we have to define a pattern functor [15] for expressions

```
\mathbf{data} \; \mathsf{Expr}_\mathsf{F} \; \mathsf{r} = \mathsf{Lit}_\mathsf{F} \; \mathsf{Int} \; | \; \mathsf{Add}_\mathsf{F} \; \mathsf{r} \; \mathsf{r}
```

that has the same structure as Expr_D, but abstracts from recursive calls. We furthermore have to instantiate a class MuRef to make the relationship between Expr and Expr_E precise.

Using the function reifyGraph we can convert a value of type Expr_D into a conventional graph representation based on a list of type [(Int, Expr_F Int)] associating integer labels with partial terms.

For example, reifyGraph (tree, 1) returns the graph

Graph $[(1,Add_F 2 2),(2,Lit_F 1)] 1$ where the final 1 points to the root node.

We now define a function build that transforms such a list of nodes into an explicitly shared ClosedExpr:

```
\begin{array}{l} \text{build} :: [(\text{Int}, \text{Expr}_F \, \text{Int})] \to \text{Int} \to \text{ClosedExpr} \\ \text{build env root} = \\ \text{letrec}\_(\lambda \text{vs} \to \text{let go} \, (\text{Lit}_F \, \text{x}) = \text{Lit} \, \text{x} \\ \text{go} \, (\text{Add}_F \, \text{v}_1 \, \text{v}_2) = \\ \text{Add} \, (\text{var vs v}_1) \, (\text{var vs v}_2) \\ \text{in map} \, (\text{go} \circ \text{snd}) \, \text{env}) \\ (\lambda \text{vs} \to \text{var vs root}) \\ \text{where} \\ \text{var vs n} = \\ \text{fromJust} \, (\text{lookup n} \, (\text{zipWith} \, (\lambda (\text{i},\_) \, \text{x} \to (\text{i},\text{x})) \, \text{env vs})) \end{array}
```

In this definition, var associates the integer labels with a variable from the list vs, and then looks up the label n. We convert between values of type Expr_F a and Expr_D a using the function go.

Using build, we can now write programs like

```
\begin{aligned} \text{test} = \textbf{do} \; (\text{Graph env r}) \leftarrow \text{reifyGraph (tree}_{l} \; 3) \\ \text{print (text (build env r))} \end{aligned}
```

least two situations in which the type safety we are able to obtain is not satisfactory vet. Firstly, if the language itself has a type system, then we might want to have a datatype explicitly encoding welltyped terms, which has consequences on how we have to define the

have to perform a lazy pattern match.

that observe sharing (such as text).

In the following, we will show how to fix these issues by assigning more precise types to our language constructs. Typed Lists

and then convert it to a value of type ClosedExpr using reifyGraph and build. We can then process the resulting ASG with functions

Summary Our ASG representation is suitable for representing various binding constructs in Haskell DSLs. However, there are at

binding constructs. Secondly, for mutually recursive bindings we can either add on a constructor for each number of bindings and go via tuples, or we can add one constructor working with lists as we have done. However, this requires maintaining an implicit invariant that we match on no more bindings than we are defining, and we

This section presents typed lists. Typed lists are a generalization of both homogeneous and heterogeneous lists of statically known length. We will make use of typed lists for encoding well-typed mutually recursive bindings in Sections 5 and 6.

Typed lists are defined using the following datatype: data TList:: $(* \rightarrow *) \rightarrow * \rightarrow *$ where

TNil::TList f() $(:::)::ft \rightarrow TListfts \rightarrow TListf(t,ts)$ (t,ts) representing the list with t as the head and ts as the tail.² The signature determines both the length of the typed list and the types of its elements. Where the signature contains a type t, the corresponding element has type f t.

Heterogeneous and homogeneous lists Typed lists can be viewed

A typed list TList f ts is parameterized by a type constructor f of kind $* \to *$ and indexed by a *signature* of types ts. The signature encodes a type-level list, with () representing the empty list and

as a generalization of heterogeneous lists of statically known length. Heterogeneous lists correspond to the case where f=I, and I is the identity type constructor:

 $\mathsf{hlist} :: \mathsf{TList} \ \mathsf{I} \ (\mathsf{Int}, (\mathsf{Int} \to \mathsf{Int}, (\mathsf{Bool}, ())))$

Using I we can encode the following heterogeneous list:

newtype $| a = | \{unl :: a\}$

hlist = $13 ::: 1 (\lambda x \rightarrow x) ::: 1$ False ::: TNil

In this case hlist is an heterogeneous list that contains values of type

Int, Int → Int and Bool as elements, and the types of the elements are reflected in the signature.

Typed lists are also a generalization of homogeneous lists. Ho-

mogeneous lists correspond to the case where a = K b, and K b is the constant type constructor:

newtype K b a = K $\{unK :: b\}$

For example, we can encode the list [1,2,3] as follows:

² Alternatively, we could use recent GHC extensions that allow kind polymorphism and datatype *promotion* [28] to provide a more direct definition of typed lists:

Alternatively, we could use recent GHC extensions that all

```
 \begin{split} & \textbf{data} \; \text{TList} :: (k \to *) \to [k] \to * \; \textbf{where} \\ & \text{TNil} :: \; \text{TList} \; f \; '[] \\ & (:::) :: \; f \; t \to \; \text{TList} \; f \; ts \to \; \text{TList} \; f \; (t \; \; ': \; ts) \end{split}
```

```
list = K1:::K2:::K3:::TNii

The use of the constant functor means that all elements are of two lot. The concrete types that occur in the signature become
```

type Int. The concrete types that occur in the signature become irrelevant; the signature merely encodes the length of the list.

*Basic operations** We can access the head and the tail of non-

Basic operations We can access the head and the tail of nonempty typed lists: thead::TList $f(t,ts) \rightarrow ft$

thead (x ::: xs) = xttail :: TList $f(t,ts) \rightarrow TList fts$ ttail (x ::: xs) = xs

list:: TList (K Int) $(t, (t_1, (t_2, ())))$

Unlike for regular head and tail, no pattern matching errors can occur in thead and ttail, because the type signature specifies that

the input list must have at least one element.

Another useful operation is tlength, which returns the number of elements in a typed list:

tlength :: TList v t \rightarrow Int tlength TNil = 0 tlength (x:::xs) = 1 + tlength xs

Mapping and zipping Operations like map or zipWith have counterparts in the world of typed lists. Where map lifts a function of type $a \rightarrow b$ to a function on lists, the corresponding tmap operates on a *natural transformation* of type $\forall t.ft \rightarrow gt$:

tmap :: $(\forall t.f t \rightarrow g t) \rightarrow TList f t \rightarrow TList g t$ tmap h TNil = TNil tmap h (x ::: xs) = h x ::: tmap h xs

Apart from the more general type, the code of tmap is the same as that for map. We can easily obtain a specialized version for homogeneous lists:

tmapK:: $(a \rightarrow b) \rightarrow TList (K a) ts \rightarrow TList (K b) ts$ tmapK $f = tmap (K \circ f \circ unK)$

A generalization of zipWith for typed lists can be obtained in a similar fashion:

Note that the type signature of tzipWith dictates that both input lists as well as the output list share a common signature and therefore must in particular be of the same length. As a result, we have to provide only two cases, where either both input lists are empty, or both input lists are non-empty.

Producers of typed lists We will also need a version of iterate that operates on typed lists. This operation is interesting because it *produces* a typed list, whereas all the functions we have defined above are *consumers* of typed lists.

While the conventional iterate function produces an infinite list, we now have to produce a list of a statically given signature, and in particular length. We therefore have to define our typed version of iterate by induction over the signature ts. As a consequence, the function cannot simply be of type

$$(a \rightarrow a) \rightarrow a \rightarrow TList (K a) ts$$

because we have to produce a result that is polymorphic in ts, and we have no way in Haskell to analyze ts. We can, however, use a well-known type-level programming technique [6] to reflect the structure of the signature to the value level and then perform induction over the reflected signature:

data RList:: $* \rightarrow *$ where

RNil :: RList ()

RCons:: RList ts \rightarrow RList (t, ts)

Using RList, it is now straight-forward to define a version of iterate for typed lists:

titerate' :: RList ts \rightarrow (a \rightarrow a) \rightarrow a \rightarrow TList (K a) ts titerate' RNil f n = TNil titerate' (RCons xs) f n = K n ::: titerate' xs f (f n)

Using type classes for producers Using titerate' is inconvenient, because in order to invoke it, we have to pass a term of type RList ts, and constructing such a term is tedious. We can, however, use a type class to build a value of the appropriate type automatically and

pass it implicitly, so that we can define a more convenient function

titerate :: CList ts \Rightarrow (a \rightarrow a) \rightarrow a \rightarrow TList (K a) ts titerate = titerate' cList

The type class CList and its instances are:

class CList t where

titerate as follows:

instance CList () where clist = RNil

instance CList ts \Rightarrow CList (t,ts) where cList = RCons cList

The resulting function titerate can be used almost in the same way as iterate:

```
tenumFrom :: CList ts \Rightarrow Int \rightarrow TList (K Int) ts tenumFrom n = titerate (+ 1) n test :: TList (K Int) (t_1,(t_2,())) test = tenumFrom 0
```

The main difference is that the type is important to determine how many elements will be generated. For example, test generates a list with the elements $K\ 0$ and $K\ 1$, because the signature of test is a type-level list with two elements t_1 and t_2 .

5. Typed ASGs and DSLs

We will illustrate this by adapting the interpreter presented throughout Section 3 to ensure that all terms are well-typed by construction. As for the untyped interpreter, observing sharing and recursion is possible. Because mutually recursive LetRec subsumes normal Let and Mu, we drop the latter two from the language.

This section shows how to define well-typed abstract syntax graphs.

Well-typed Abstract Syntax Graphs If we want to model well-typed ASGs, we have to first introduce an additional type argument that serves as the index for the type of the value being represented, and then adapt the types of the constructors in order to establish the typing rules of the embedded language.

But how do we represent variables? As the embedded language is now indexed by a type argument, variables can be of different (Haskell) types. Therefore, we change the type parameter for variables from kind * to kind * \rightarrow *: we pass in a type function that, given a type of the embedded language, returns the associated type of variables. If we apply this strategy to our example expression language, we end up with the following datatype:

```
\label{eq:type-closed-expr} \begin{split} & \text{type-closed-expr} \ t = \forall f. \text{Exprft} \\ & \text{data Expr} \ (f ::: * \to *) :: * \to * \text{ where} \\ & \text{Lit} \qquad :: \text{Int} \to \text{ExprfInt} \\ & \text{Add} \qquad :: \text{ExprfInt} \to \text{ExprfInt} \to \text{ExprfInt} \\ & \text{IfZero} \quad :: \text{ExprfInt} \to \text{Exprft} \to \text{Exprft} \\ & \text{Var} \qquad :: \text{ft} \to \text{Exprft} \\ & \text{Lam} \qquad :: (ft_1 \to \text{Exprft}_2) \to \text{Exprf} \ (t_1 \to t_2) \\ & \text{App} \qquad :: \text{Exprf} \ (t_1 \to t_2) \to \text{Exprft}_1 \to \text{Exprft}_2 \\ & \text{LetRec} :: \text{CList ts} \Rightarrow (\text{TList fts} \to \text{TList} \ (\text{Exprf}) \ \text{ts}) \to \\ & \qquad \qquad (\text{TList fts} \to \text{Exprft}) \to \text{Exprft} \end{split}
```

In the Var case, we pass the type t to the parameter function f to obtain a suitable variable type, as was our plan. The case for Lam shows that apart from adding type arguments everywhere, the structure of representing binders remains the same.

The case for mutually recursive bindings LetRec is more interesting. We now use typed lists (as introduced in Section 4) rather than ordinary lists. They keep track of the types of all the elements in the list, and at the same time determine the length of the list. Therefore, by using the same signature ts three times for the three occurrences of TList, we now establish *statically* that all three occurrences have exactly the same shape. This is a big improvement over the untyped encoding which does not provide such guarantees.

Furthermore, the CList ts constraint in LetRec guarantees that expressions built with this constructor support reifying the typelevel list into a value of type RList ts. This is useful when we want to use producer functions like titerate to define functions over Expr.

As in the untyped setting parametricity still ensures that we

As in the untyped setting, parametricity still ensures that we cannot inspect variables as long as an expression is polymorphic in the variable type function f. We define ClosedExpr as an abbreviation for such closed terms again.

Well-typed evaluator The code for the well-typed evaluator is:

```
\begin{array}{lll} \text{eval (Lit i)} & = \text{i} \\ \text{eval (Add } e_1 \, e_2) & = \text{eval } e_1 + \text{eval } e_2 \\ \text{eval (IfZero } e_1 \, e_2 \, e_3) = \text{if eval } e_1 == 0 \text{ then } \text{eval } e_2 \text{ else } \text{eval } e_3 \\ \text{eval (Var x)} & = \text{unl x} \\ \text{eval (Lam e)} & = \text{eval } \circ e \circ \text{I} \\ \text{eval (App } e_1 \, e_2) & = \text{(eval } e_1) \text{ (eval } e_2) \\ \text{eval (LetRec es e)} & = \text{eval (e (fix (tmap (I \circ \text{eval}) \circ \text{es})))} \end{array}
```

eval:: Expr I $t \rightarrow t$

Unlike the interpreter we defined in Section 3, there is no need for a separate Value datatype for values. Since we used the Haskell type constructors Int and (→) to model the type language of the embedded language, we can simply use t as value type of a term that has type Exprft. For the purposes of evaluation, we have to instantiate f with a type function that makes this relation explicit: the identity type constructor!

The resulting interpreter is untagged. There are no constructors

wrapping the values, and we do not need to perform any typechecking at run-time. We statically know that in each construct, the arguments we obtain are of the correct types.

While the code for the "normal" language constructs becomes simpler, the code for the binding constructs remains nearly unchanged: we only have to sprinkle coercion functions unl and I to help the type checker along.

As before, we can define wrappers for certain constructors to make the use of the language a bit more convenient. For example:

```
 \begin{split} (\oplus) &= \mathsf{Add} \\ \mathsf{one} &= \mathsf{Lit} \ 1 \\ \mathsf{lam}\_{::} \left( \mathsf{Exprft}_1 \to \mathsf{Exprft}_2 \right) \to \mathsf{Exprf} \left( \mathsf{t}_1 \to \mathsf{t}_2 \right) \\ \mathsf{lam}\_{\,e} &= \mathsf{Lam} \left( \lambda \mathsf{x} \to \mathsf{e} \left( \mathsf{Var} \, \mathsf{x} \right) \right) \\ \mathsf{letrec}\_{::} \mathsf{CList} \ \mathsf{ts} \Rightarrow \left( \mathsf{TList} \left( \mathsf{Exprf} \right) \ \mathsf{ts} \to \mathsf{TList} \left( \mathsf{Exprf} \right) \ \mathsf{ts} \right) \to \\ &\qquad \qquad \left( \mathsf{TList} \left( \mathsf{Exprf} \right) \ \mathsf{ts} \to \mathsf{Exprft} \right) \to \mathsf{Exprft} \end{aligned}
```

```
letrec_ es e = LetRec (\lambda xs \rightarrow es (tmap Var xs))
                                 (\lambda xs \rightarrow e \text{ (tmap Var } xs))
```

If we want non-recursive let-bindings or a simple fixed-point construct back, we can easily define these in terms of letrec_: let_{\perp} :: Expr f $t_1 \rightarrow (Expr f t_1 \rightarrow Expr f t_2) \rightarrow Expr f t_2$

94

```
let_e_1 e_2 = letrec_(\lambda_- \rightarrow e_1 ::: TNil) (\lambda(x ::: TNil) \rightarrow e_2 x)
Using the definitions above, we can still distinguish between im-
```

plicitly shared and explicitly shared terms as before. The two versions tree, and tree, defined in Figure 1 are valid in the typed setting without any change of the code - only the type becomes

Int \rightarrow ClosedExpr Int

in both cases. The implicitly shared version will still lose sharing, whereas the explicitly shared version still evaluates quickly.

and recursion apply to well-typed terms: the addition of typing information does not affect the preservation of sharing and recursion. On the other hand, we now can no longer define terms that are ill-typed according to the type system of our DSL. For instance,

In summary, the same properties regarding observable sharing

example from Section 3 fails to type check, because it uses twice at two different types, but our DSL has only monomorphic types.

Printing terms Let us also look at how we have to adapt the function text that we have introduced in Section 2. Here, the relation between DSL types and result types is different compared to evaluation: regardless of the DSL type that a variable has, they are all printed as strings. We need the K type constructor instead of I: text::ClosedExpr $t \rightarrow String$

```
go :: Expr (K String) t \rightarrow Int \rightarrow String
               -- cases for Lit, Add, IfZero, App as before
         go (Var x)
                     _{-} = unK x
        go (Lam e) c =
           "(\\" ++ v ++ " -> " ++ qo (e (K v)) (c + 1)
            where v = "v" ++ show c
        go (LetRec es e) c =
            "(let { " ++ intercalate "; " ds ++
           " } in "++go (e vs) c'++")"
            where
              vs = tmapK (\lambda i \rightarrow "v" ++ show i) (tenumFrom c)
              c' = c + t length vs
              ds = ttoList$
                    tzipWith (\lambda(Kv) e \rightarrow K(v ++ " = " ++ go e c'))
                            vs (es vs)
Similarly to the evaluator code, we must add a few coercion func-
```

text e = go e 0

tions (K and unK) throughout the pretty printer code.

The LetRec case is interesting again, because we have to deal with typed lists. In vs, we define the strings representing each of the bound variables. First, we generate numbers starting from

the current counter c using tenumFrom. Then we map over the list, moving from type K Int to K String. How many variables are

generated is determined by the type context! In the declaration of ds, we pass our typed list of strings to the declaration function es, and the type of LetRec dictates that the input lists of variables must have the same shape as the output list of bindings.

Note that tenumFrom works only for result types that actually are list types, as witnessed by the CList constraint – this is an

In ds, we then take the list of variables vs and the list of expressions es vs and generate strings representing each of the bindings using tzipWith. We end up with a typed list containing elements of type K String, but we would actually like to have a list of strings at this point. The function ttoList achieves this:

example for why we need to put a CList constraint in the type of

the LetRec constructor.

ttoList::TList (K a) ts \rightarrow [a] ttoList TNil = [] ttoList (K x:::xs) = x:ttoList xs Finally, we separate each of the bindings by "; " by using the standard list function intercalate and append everything together in a single string.

6. Encodings of ASGs

In this section, we discuss an encoding of ASGs using type classes. This approach is interesting because it stands somewhere inbetween a *shallow* and a *deep* embedding. Like for deep embeddings, it is possible to have multiple interpretations and perform a form of syntactic analysis. Like for shallow embeddings, it is possible to to create (but not observe) sharing using Haskell's built-in let. Moreover, it is easy to extend the language and add new con-

structs without touching existing code. For the deep embeddings we have been using in Sections 2, 3 and 5, adding a new construc-

tor requires modifying all the interpretation functions.

Uses of sharing in the embedded language can still be explicit and therefore can be observed and as needed and we can maintain the level of type safety established in Section 5.

An additional advantage of the class-based encoding is that it is possible to provide reusable code for the binding constructs we have presented. As binding constructs are useful and similar

throughout many DSLs, being able to reuse code reduces the implementation burden on DSL designers.

Encoding datatypes as type classes Hinze [12] showed that type classes provide a convenient way to define Church encodings of datatypes. Moving from a (generalized) algebraic datatype to a class is an entirely mechanical process [21]. As an example, let us consider how to encode simple well-typed arithmetic expressions like the ones presented in Section 5, i.e., we base this construction on the type Expr from Section 5, but we consider only the Lit and

Add constructors for now:

class ArithAlg expr $(f::* \rightarrow *)$ where

lit :: Int \rightarrow expr f Int

 $(\oplus) :: \mathsf{expr}\,\mathsf{f}\,\mathsf{Int} \to \mathsf{expr}\,\mathsf{f}\,\mathsf{Int} \to \mathsf{expr}\,\mathsf{f}\,\mathsf{Int}$

Looking at the transformation from a syntactic perspective, all we did was: change the datatype declaration to a class, transform the data constructors into methods of the class, and use a class parameter expr wherever the original datatype Expr was being used.

Semantically, ArithAlg encodes the signature of algebras of the original datatype. Instances of ArithAlg correspond to fold-like functions over that datatype. For the reader interested in knowing more about this technique we suggest several resources available elsewhere [5, 12, 21].

Encoding Binders We can follow the same recipe to encode binders. Let us again consider the datatype Expr from Section 5, ignoring all but the two binding-related constructors Var and LetRec. We obtain the following type class:

```
class BindAlg expr f where
```

```
var :: ft \rightarrow exprft
letrec :: CList ts \Rightarrow (TList fts \rightarrow TList (exprf) ts) \rightarrow
(TList fts \rightarrow exprft) \rightarrow exprft
```

As before, it is possible to define wrappers that allow more convenient use of binding constructs, or that define simpler binding constructs in terms of letree. As an example, here is the code for non-recursive let-bindings:

```
\begin{split} \text{let}\_{::} \, \text{BindAlg expr} \, f &\Rightarrow \text{expr} \, f \, t_1 \to \\ & (\text{expr} \, f \, t_1 \to \text{expr} \, f \, t_2) \to \text{expr} \, f \, t_2 \\ \text{let}\_{\, e_1} \, e_2 &= \text{letrec} \, (\lambda_- \to e_1 \, \text{:::} \, \text{TNil}) \, (\lambda(x \, \text{:::} \, \text{TNil}) \to e_2 \, (\text{var} \, x)) \end{split}
```

this generic behavior by providing a "default" instance for BindAlg. This instance can be reused when defining suitable operations: $\textbf{newtype} \ \text{Default} \ f \ t = D \ \big\{ \text{unD} :: f \ t \big\}$

Generic behavior for binders As observed by Oliveira and Cook [20], there are many operations that share a common definition for the binding constructs. We use this observation to capture

instance BindAlg Default f where var x = D x letrec es e = e (fix (tmap unD ∘ es))

letrec es e = e (lix (timap unidoes)

This definition turns out to be useful for operations such as evaluation or inlining – whenever we do not need to observe sharing. For functions such as text, we want to observe the binding structure and will require a different instance.

will require a different instance.

Extensibility and Modularity As Oliveira et al. [22] shows, an advantage of using the class-based approach is that in contrast to datatypes, which are closed to extension, we can add new cases to

a language simply by defining another class. In other words this technique can be used to solve the *expression problem* [27].

Several DSLs share a number of common components. For example, many DSLs will have the arithmetic expressions and recur-

sive binders that we already discussed. Adding a new set of DSL constructs is as simple as defining a new type class. For example, we can create a third class LamAlg for lambda and application:

class LamAlg expr f where

lam:: $(ft_1 \rightarrow exprft_2) \rightarrow exprf(t_1 \rightarrow t_2)$ app:: exprf $(t_1 \rightarrow t_2) \rightarrow exprft_1 \rightarrow exprft_2$

If we want to state that expressions are given by the combination of the three classes we have defined above, we can denote that with

```
ExprAla expr f
We can build terms in this language by applying and combining
class methods from each of the three classes. Closed expressions
```

class (BindAlg expr f, ArithAlg expr f, LamAlg expr f) ⇒

type ClosedExpr $t = \forall expr f. ExprAlg expr <math>f \Rightarrow expr f t$

are overloaded in the instantiation of ExprAla:

Evaluation In order to define evaluation, we define instances for each of the classes separately. If desired, we could define these instances at different times, when we decide to extend the language with new constructs, and without touching existing code.

No instance for BindAlg is needed - we can reuse the default instance defined above. We have to provide instances for ArithAlg

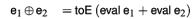
and LamAlg, however. A ClosedExpr t evaluates to a t, so we could choose I as the instantiation of the f arguments of the classes. However, while several functions on expressions might share the

same type signature, there can be only a single instance each for ArithAlg Default I and LamAlg Default I. Therefore, we define a new type isomorphic to I specifically for the evaluation function:

```
eval :: Default Eval t \rightarrow t
eval = unF \circ unD
toE ::t → Default Eval t
toE = D \circ E
```

The instances are then straightforward: instance ArithAlg Default Eval where lit i = to E i

newtype Eval $t = E \{unE :: t\}$



```
= toE (eval \circ f \circ E)
      app e_1 e_2 = toE ((eval e_1) (eval e_2))
    instance ExprAlg Default Eval
Note that eval can be applied directly to a term of type ClosedExpr.
Shared trees If we abbreviate one = lit 1, and use Int \rightarrow ClosedExpr Int
as the type signature, then the two versions tree, and tree from
Figure 1 work once again. However, the behavior here is similar
to what we discussed in Section 2 for shallow DSLs; both versions
preserve sharing. There is no indirection of data constructors when
using the class-based encoding: a term is directly encoded as its
interpretation (or actually, all possible interpretations).
   Still, there are advantages to using explicit sharing, as implicit
sharing remains rather fragile: if we, for example, define a function
    double :: ClosedExpr t → ClosedExpr t
that traverses an expression and doubles all literals, then the traver-
sal over an implicitly shared term will destroy the sharing. We also
need to be explicit whenever we have to observe sharing.
Printing terms As an example of an operation that observes shar-
ing and does not make use of the default instance for BindAlq, we
return to our text function. We only show the instance for BindAlg:
    newtype Text (f::* \rightarrow *) t = T \{ text' :: Int \rightarrow String \}
    instance BindAlg Text (K String) where
      \operatorname{var} x = T(\lambda_{-} \to \operatorname{unK} x)
      letrec es e = T (\lambda c \rightarrow
         let vs = tmapK (\lambda i \rightarrow "v" ++ \text{show i}) (tenumFrom c)
            c' = c + t length vs
            ds = ttoList$
                   tzipWith (\lambda(K v) e \rightarrow K (v ++ " = " ++ text' e c'))
                            vs (es vs)
         in "(let { " ++ intercalate "; " ds ++
```

instance LamAlg Default Eval where

```
" } in " ++ \text{text'} (e \text{ vs}) c' ++ ")")
```

The code is nearly the same as that given in Section 5. The actual text function wraps text':

 $\begin{array}{l} \text{text} :: \forall t. \text{ClosedExpr} \: t \rightarrow \text{String} \\ \text{text} \: e = \text{text'} \: (e :: \text{Text} \: (\text{K String}) \: t) \: 0 \end{array}$

Summary Using the class-based encoding of ASGs is recommended whenever extensibility is a must. It is easily possible to encode well-typed terms in the class-based setting, but actually not necessary. The untyped ASGs of Section 3 can easily be translated into the class-based setting as well. A disadvantage of the class-based encoding is that it forces interpretation functions to be folds – if functions require nested pattern matches or have strange recursion behavior, they can be tricky to encode as algebras.

7. Conclusion

We propose using ASGs instead of ASTs to provide a more convenient representation of abstract syntax for EDSLs. ASGs internalize the information about sharing and recursion directly in the representation. As such, environments can in most cases be avoided, and there is no need to deal with other binding-related issues such as α -equivalence, name capture or generation of fresh variables.

We show that ASGs are flexible: they extend nicely to the generalized setting of terms with statically encoded type information, and they are suitable for class-based as well as datatype-based encodings. The ability to encode well-typed terms is particularly interesting, because many DSLs have type systems that can be encoded directly in the host language. The ability to modularize DSL constructs is important to provide flexibility and reuse.

References

- [1] T. Altenkirch and B. Reus. Monadic presentations of lambda terms using generalized inductive types. In CSL '99, 1999. [2] A. I. Baars and S. D. Swierstra. Type-safe, self inspecting code. In
- Haskell '04, 2004. [3] A. I. Baars, S. D. Swierstra, and M. Viera. Typed transformations of typed grammars: The left corner transform. Electron. Notes Theor.
- Comput. Sci., 253(7), 2010. [4] P. Biesse, K. Claessen, M. Sheeran, and S. Singh. Lava: hardware design in Haskell. In ICFP '98, 1998.
- [5] J. Carette, O. Kiselyov, and C. Shan. Finally tagless, partially evaluated: Tagless staged interpreters for simpler typed languages. J. Funct. Program., 19(5), 2009.
- [6] J. Cheney and R. Hinze. A lightweight implementation of generics and dynamics. In Haskell 2002, 2002.
- [7] A. Chlipala. Parametric higher-order abstract syntax for mechanized semantics. In ICFP'08, 2008. [8] K. Claessen and D. Sands. Observable sharing for functional circuit description. In In Asian Computing Science Conference, Springer
- Verlag, 1999. [9] D. Devriese and F. Piessens. Finally tagless observable recursion for an abstract grammar model. Journal of Functional Programming, 22(6):757-796, November 2012.
- [10] D. Devriese, I. Sergey, D. Clarke, and F. Piessens. Fixing idioms: A recursion primitive for applicative dsls. In PEPM'13, 2013. [11] A. Gill. Type-safe observable sharing in Haskell. In *Haskell'09*, 2009.
- [12] R. Hinze. Generics for the masses. J. Funct. Program., 16(4-5), 2006.
- [13] C. Hofer, K. Ostermann, T. Rendel, and A. Moors. Polymorphic

[15] P. Jansson and J. Jeuring. Polyp – a polytypic programming language extension. In *POPL'97*, 1997.
[16] O. Kiselyov. Implementing explicit and finding implicit sharing in embedded DSLs. In *Proceedings IFIP Working Conference on Domain-*

[14] P. Hudak. Building domain-specific embedded languages.

embedding of dsls. In GPCE '08, 2008.

Computing Surveys, 28, 1996.

Specific Languages, 2011.

- [17] E. Meijer and G. Hutton. Bananas in space: extending fold and unfold to exponential types. In FPCA'95, 1995.
 [18] M. Might, D. Darais, and D. Spiewak. Parsing with derivatives: a functional pearl. In ICFP '11, 2011.
 [19] J. T. O'Donnell. Overview of Hydra: A concurrent language for
- synchronous digital circuit design. In *IPDPS '02*, 2002.
 [20] B. C. d. S. Oliveira and W. R. Cook. Functional programming with structured graphs. In *ICFP '12*, 2012.
 [21] B. C. d. S. Oliveira and J. Gibbons. Typecase: a design pattern for type-indexed functions. In *Haskell '05*, 2005.
- [21] B. C. d. S. Oliveira and J. Gibbons. Typecase: a design pattern for type-indexed functions. In *Haskell '05*, 2005.
 [22] B. C. d. S. Oliveira, R. Hinze, and Andres Löh. Extensible and Modular Generics for the Masses. In *TFP '06*, 2006.
 [23] S. Peyton Jones, D. Vytiniotis, S. Weirich, and G. Washburn. Simple
- [23] S. Peyton Jones, D. Vytiniotis, S. Weirich, and G. Washburn. Simple unification-based type inference for gadts. In *ICFP'06*, 2006.
 [24] F. Pfenning and C. Elliot. Higher-order abstract syntax. In *PLDI '88*, 1988.
- [25] F. Pottier. Lazy least fixed points in ML. Unpublished, 2009.[26] P. Wadler. The essence of functional programming. In *POPL'92*, 1992.
- [27] P. Wadler. The essence of functional programming. In *POPL'92*, 1992.
- [27] I. Wallet. The expression problem. Note to sava Genericky maining list, Nov. 1998.
 [28] B. A. Yorgey, S. Weirich, J. Cretin, S. Peyton Jones, D. Vytiniotis, and J. P. Magalhães. Giving Haskell a promotion. In *TLDI '12*, 2012.