

Trees that Evolve Together

Shayan Najd

Laboratory for Foundations of Computer Science
The University of Edinburgh, Scotland, U.K.
sh.najd@gmail.com

August 8, 2018

Abstract

We study the notion of evolution in functional data types, as a new approach to the problem of decorating abstract syntax trees with additional information. We observed the need for such evolutions while redesigning the data types representing Haskell abstract syntax inside Glasgow Haskell Compiler (GHC).

Specifically, we generalise the Trees that Grow (TTG) idiom, to account for more refined forms of evolution in abstract syntax trees. The generalisation, while introduces a modest change, results in a remarkably more powerful approach. For example, the new approach can account for changes in the type of specific, but not all, subterms.

We show how such newly gained powers help with solving some practically challenging problems, such as handling source location decorations, introducing mutual recursion, and lightweight generic programming.

1 Introduction

A forest with one tree, is no forest; there are many trees in a forest, all growing together in harmony by passing of time.

A compiler with one abstract syntax data type, is no compiler; there are many abstract syntax data types in a compiler, all evolving together synchronously through passes of the compiler.

Trees that Grow (TTG) idiom (Najd and Peyton Jones, 2017) introduces a practical solution to the problem of extending abstract syntax data types in a functional setting. TTG supports a wide range of extensions to normal algebraic data types, and scales to support generalised algebraic data types. Extensions can account for the evolution of an abstract syntax data type through compiler passes, by assigning a different set of extensions for each pass-specific variant of an abstract syntax data type. For a specific pass, these sets of extensions can be synchronised across data types by sharing the same parameter used for describing the set of extensions. The same synchronisation process can be used to independently extend mutually-defined data types, straightforwardly.

Less straightforward is when two independently defined data types evolve to be mutually-defined.

For example, in Glasgow Haskell Compiler (GHC), expressions in the parsing, or type checking phase, are decorated with source locations, but the same form of expressions are entirely undecorated when used in meta programming, i.e, Template Haskell. Consider the following simplified data types describing expression in application form, in the parsing pass (denoted by $-^P$), type checking pass (denoted by $-^T$), and meta programming (denoted by $-^M$).

<pre> type $LExp^P$ $id = Loc (Exp^P id)$ data $Exp^P id$ = $App^P (LExp^P id) (LExp^P id)$ \vdots </pre>	<pre> type $LExp^T id = Loc (Exp^T id)$ data $Exp^T id$ = $App^T Typ (LExp^T id) (LExp^T id)$ \vdots </pre>
<pre> data $Exp^M id$ = $App^M (Exp^M id) (Exp^M id)$ \vdots </pre>	

where the wrapper for source location spans is defined as follows.

```

type  $Spn = (Integer, Integer)$       | data  $Loc trm = Loc Spn trm$ 

```

We really would like to have, instead of three data types, one single base data type describing expressions, and define these three trees as instances of the base data type. This is exactly the goal of the TTG idiom. Following the TTG idiom, the base data type can be defined as the following.

<pre> type $LExp^X x = Loc (Exp^X x)$ data $Exp^X x$ = $App^X !(X_{App} x) (LExp^X x) (LExp^X x)$ $New^X !(X_{New} x)$ \vdots </pre>	<pre> class $X_{Exp} x$ where type $X_{App} x$ type $X_{New} x$ \vdots </pre>
---	---

The base data type Exp^X clearly subsumes Exp^P and Exp^T by instantiating the extension type families such as X_{App} :

<pre> data P id type $LExp^P$ $id = LExp^X (P\ id)$ type Exp^P $id = Exp^X (P\ id)$ instance $X_{Exp} (P\ id)$ where type $X_{App} (P\ id) = NoFld$ type $X_{New} (P\ id) = NoCon$ ⋮ </pre>	<pre> data T id type $LExp^T$ $id = LExp^X (T\ id)$ type Exp^T $id = Exp^X (T\ id)$ instance $X_{Exp} (T\ id)$ where type $X_{App} (T\ id) = Typ$ type $X_{New} (T\ id) = NoCon$ ⋮ </pre>
---	---

where the empty field extension type and the empty constructor extension type are defined as the following:

<pre> data $NoFld = NoFld$ </pre>	<pre> data $NoCon$ </pre>
---	---

It is straightforward to provide pattern synonyms, as the following, to hide the plumbing, which in turn makes the one-to-one correspondence evident.

```

pattern  $App^P :: LExp^P\ id \rightarrow LExp^P\ id \rightarrow Exp^P\ id$ 
pattern  $App^P\ l\ m = App^X\ NoFld\ l\ m$ 
⋮
pattern  $App^T :: Typ \rightarrow LExp^T\ id \rightarrow LExp^T\ id \rightarrow Exp^T\ id$ 
pattern  $App^T\ a\ l\ m = App^X\ a\ l\ m$ 
⋮

```

Nevertheless, the process is not as straightforward when it comes to subsuming Exp^M by instantiating Exp^X . The subexpression in Exp^X are of type $LExp^X$; the source location decorations are hard coded inside the subexpressions of Exp^X , deeming Exp^X an invalid choice for a base data type for subsuming Exp^M . Can we find a suitable base data type that subsumes all the three data types at once? Such a base data type should also be able to subsume data types reasonably similar to these three. What are the general guidelines?

In this paper, we try to answer these questions, by generalising the TTG idiom, to support changes to the type of specific, not all, subexpressions. This general case, handles subsumption of Exp^P , Exp^T and Exp^M by one single base data type, where the non-trivial change (beyond the simpler forms already supported by TTG) is the evolution of the type of some subexpressions from Exp^X to $LExp^X$. For example, in the App^X constructor of the ideal base type both arguments change type from Exp^X to $LExp^X$ when instantiated for Exp^P and Exp^T , but remain untouched for Exp^M . Such ability to change the type of specific subterms is beyond the original TTG approach, but as we show is easy achievable by a modest generalisation.

To see the necessity of such generalisation, let us attempt to solve the problem of subsuming at once all the three data types Exp^P , Exp^T , and Exp^M by an approach dual to the one taken earlier: we start with a base data type where

source location decorations are not hard coded into subexpressions.

Following the original TTG idiom, we have the following fully undecorated base data type:

<pre> data Exp^X x = App^X $!(X_{App} \ x) \ (Exp^X \ x) \ (Exp^X \ x)$ New^X $!(X_{New} \ x)$ ⋮ </pre>	<pre> class X_{Exp} x where type X_{App} x type X_{New} x ⋮ </pre>
--	--

Now to store the source location decorations for the subexpressions in Exp^P and Exp^T , we have two options: either we store them in the extension fields, or in the extension constructor. The former leads to a new programming interface that barely corresponds to the original one. In the former design, the distinction between undecorated and decorated subexpressions is destroyed, leading to unnecessary couplings of code handling decorations and undecorated expressions. For example, with the former design, source location decorations should be handled in every clause of matchings, whenever we pattern match on an expression; it unnecessarily couples the task of handling source location decorations with any other task inspecting expressions. The latter design, is more promising in that handling source location decorations remains decoupled from other task in that when matching on expressions only one clause needs to handle source location decorations, i.e., the clause handling the constructor extension. However, as we discuss shortly, the latter design has its own drawbacks.

The following is how we instantiate the base data type Exp^X , for the three data type, following the latter design of storing source location decorations in the constructor extension.

<pre> data $P \ id$ type $LExp^P \ id = Exp^X \ (P \ id)$ type $Exp^P \ id = Exp^X \ (P \ id)$ instance $X_{Exp} \ (P \ id)$ where type $X_{App} \ (P \ id) = NoFld$ type $X_{New} \ (P \ id) = Loc \ (Exp^P \ id)$ ⋮ </pre>	<pre> data $T \ id$ type $LExp^T \ id = Exp^X \ (T \ id)$ type $Exp^T \ id = Exp^X \ (T \ id)$ instance $X_{Exp} \ (T \ id)$ where type $X_{App} \ (T \ id) = Typ$ type $X_{New} \ (T \ id) = Loc \ (Exp^T \ id)$ ⋮ </pre>
---	---


```

data  $M \ id$ 
type  $Exp^M \ id = Exp^X \ (M \ id)$ 
instance  $X_{Exp} \ (M \ id)$  where
  type  $X_{App} \ (M \ id) = NoFld$ 
  type  $X_{New} \ (M \ id) = NoCon$ 
  ⋮

```

We can provide pattern synonyms for constructors such App as before. Moreover, we can also provide synonyms for the wrapper at each type:

```

pattern  $Loc^P :: Spn \rightarrow Exp^P id \rightarrow LExp^P id$ 
pattern  $Loc^P sp m = New^X (Loc sp m)$ 
pattern  $Loc^T :: Spn \rightarrow Exp^T id \rightarrow LExp^T id$ 
pattern  $Loc^T sp m = New^X (Loc sp m)$ 

```

Unfortunately, above does not suffice! There are two key problems with it.

The above TTG design allows for wrapping of an already-wrapped expression, an unintended behavior that the original design (i.e., in the original data types Exp^P and Exp^T) avoids via the mutual recursion. Above design, leads to redundant, and potentially incoherent, source location decorations. For example, the following is wrongly accepted:

```

badWrap ::  $LExp^P id \rightarrow LExp^P id$ 
badWrap m =  $Loc^P (0, 0) m$ 

```

Furthermore, the above TTG design, suffers from a dual problem: it allows for expressions without source location decorations, an unintended behavior that the original design also avoids via the mutual recursion. For example, the following is wrongly accepted:

```

badSelfApp ::  $Exp^P id \rightarrow Exp^P id$ 
badSelfApp m =  $App^P m m$ 

```

What we really need is the ability to change the type of subexpressions of App^X from $Exp^X x$ to $Loc (Exp^X x')$. We can do so with a modest generalisation of the original TTG: we allow the extension parameter of subterms to be set locally by the extension type family of the constructor containing the subterms. In the original TTG idiom, extension parameters were globally set by the data type.

Here is an example of how Exp^X looks like in such a generalised setting:

<pre> data $Exp^X x$ = $\forall x_0 x_1.$ $App^X !(X_{App} x x_0 x_1) (Exp^X x_0) (Exp^X x_1)$ $New^X !(X_{New} x)$ \vdots </pre>	<pre> class $X_{Exp} x$ where type $X_{App} x x_0 x_1$ type $X_{New} x$ \vdots </pre>
--	---

Notice the only difference from earlier is that the extension parameters of subexpressions are set locally via extension type families rather than globally. For example, notice that the parameter x_0 of the first subexpression in App^X is declared locally, apart from the global x .

This modest change, has a remarkable impact on increasing the range of changes that a base data type can support. For example, here is how the new Exp^X can subsume all the three data types.

<pre> data <i>LP id</i> type <i>LExp^P id</i> = <i>Exp^X (LP id)</i> instance <i>X_{Exp} (LP id)</i> where type <i>X_{App} (LP id) x₀ x₁</i> = <i>NoCon</i> type <i>X_{New} (LP id)</i> = <i>Loc (Exp^P id)</i> ⋮ data <i>P id</i> type <i>Exp^P id</i> = <i>Exp^X (P id)</i> instance <i>X_{Exp} (P id)</i> where type <i>X_{App} (P id) x₀ x₁</i> = <i>Prf (x₀ ~ LP id</i> <i>, x₁ ~ LP id)</i> type <i>X_{New} (P id)</i> = <i>NoCon</i> ⋮ </pre>	<pre> data <i>LT id</i> type <i>LExp^T id</i> = <i>Exp^X (LT id)</i> instance <i>X_{Exp} (LT id)</i> where type <i>X_{App} (LT id) x₀ x₁</i> = <i>NoCon</i> type <i>X_{New} (LT id)</i> = <i>Loc (Exp^T id)</i> ⋮ data <i>T id</i> type <i>Exp^T id</i> = <i>Exp^X (T id)</i> instance <i>X_{Exp} (T id)</i> where type <i>X_{App} (T id) x₀ x₁</i> = <i>(Typ , Prf (x₀ ~ LT id</i> <i>, x₁ ~ LT id))</i> type <i>X_{New} (T id)</i> = <i>NoCon</i> ⋮ </pre>
--	--

```

data M id
type ExpM id = ExpX (M id)
instance XExp (M id) where
  type XApp (M id) x0 x1 = Prf ( x0 ~ M id , x1 ~ M id )
  type XNew (M id) = NoCon
  ⋮

```

We can provide pattern synonyms as before, in particular for *App* constructors, we have the following.

```

pattern AppP :: LExpP id → LExpP id → ExpP id
pattern AppP l m = AppX Prf l m
pattern AppT :: Typ → LExpT id → LExpT id → ExpT id
pattern AppT a l m = AppX (a , Prf) l m
pattern AppM :: ExpM id → ExpM id → ExpM id
pattern AppM l m = AppX Prf l m

```

where we use the following to wrap explicit proofs of equality (and other constraints in general).

```

data Prf p where
  Prf :: p ⇒ Prf p

```

In above, we define two variants of *Exp^X* for each data type *Exp^P* and *Exp^T*: one decorated with source locations at the top-level, another undecorated at the top-level. Subexpressions in one variant refers to the other alternatively, e.g., the subexpressions of *Exp^P* are of type *LExp^P* while the subexpression of *LExp^P*

is of type Exp^P . The switching between the two variants is guaranteed by the proof instance Prf in the field extension of App^X . Now, the following is rightly rejected

```
badWrap :: LExpP id → LExpP id
badWrap m = LocP (0, 0) m
```

with the following error message

```
. Couldn't match type 'LP' with 'P'
  Expected type: ExpP id
  Actual type: LExpP id
```

and, the following is also rightly rejected

```
badSelfApp :: ExpP id → ExpP id
badSelfApp m = AppP m m
```

with the following error message

```
. Couldn't match type 'P' with 'LP'
  Expected type: LExpP id
  Actual type: ExpP id
```

This pattern mimics the one-to-one correspondence between mutually recursive definition of data types and their equivalent GADT definition using indices. For example, the following mutually-recursive definition

<pre>data A = C B D A</pre>		<pre>data B = E A</pre>
-------------------------------	--	-------------------------

is in a one-to-one correspondence to the following GADT flat definition.

```
data S = A | B
data El (s :: S)
= (s ~ A) ⇒ C (El B)
| (s ~ A) ⇒ D (El A)
| (s ~ B) ⇒ E (El A)
```

But, the scope of the introduced generalisation of TTG is not limited to mutually-recursive definitions. We will explain the idea further throughout the paper.

References

Najd, S. and Peyton Jones, S. (2017). Trees that grow. *J. UCS*, 23(1):42–62.