Diffing Mutually Recursive Types

Victor Cacciari Miraldo University of Utrecht

December 8, 2016

1 Our Universe

The universe we are using is a Sums-of-Products over type variables and constant types.

```
\begin{array}{ccc} \mathsf{data} \ \mathsf{Atom} \ (n : \mathbb{N}) : \mathsf{Set} \ \mathsf{where} \\ \mathsf{I} : \mathsf{Fin} \ n & \to \mathsf{Atom} \ n \\ \mathsf{K} : \mathsf{Fin} \ ks\# & \to \mathsf{Atom} \ n \end{array}
```

Constructor I refers to the n-th type variable whereas K refers to a constant type. Value ks# is passed as a module parameter. We denote products by π and sums by σ , but they are just lists.

```
\begin{array}{l} \pi: \mathbb{N} \to \mathsf{Set} \\ \pi = \mathsf{List} \circ \mathsf{Atom} \\ \\ \sigma\pi: \mathbb{N} \to \mathsf{Set} \\ \\ \sigma\pi = \mathsf{List} \circ \pi \end{array}
```

Interpreting these codes is very simple. Here, Parms is a valuation for the type variables.

Note that here, Parms n really is isomorphic to n types that serve as the parameters to the functor $[\![F]\!]$. When we introduce a fixpoint combinator, these parameters are used to to tie the recursion knot, just like a simple fixpoint: $\mu F \equiv F(\mu F)$. In fact, a mutually recursive family can be easily encoded in this setting. All we need is n types that refer to n type-variables each!

This universe is enough to model Context-Free grammars, and hence, provides the basic bare bones for diffing elements of an arbitrary programming language. In the

future, it could be interesting to see what kind of diffing functionality indexed functors could provide, as these could have scoping rules and other advanced features built into them.

1.1 SoP peculiarities

One slightly cumbersome problem we have to circumvent is that the codes for type variables and constant types have a different *type* than the codes for types. This requires more discipline to organize our code. Nevertheless, we may wish to see Atoms as a trivial Sum-of-Product.

```
\begin{array}{l} \alpha: \{n: \mathbb{N}\} \rightarrow \mathsf{Atom} \ n \rightarrow \sigma\pi \ n \\ \alpha \ a = (a:: []):: [] \end{array}
```

Instead of having binary injections into coproducts, like we would on a regular-like universe, we have n-ary injections, or, constructors. We encapsulate the idea of constructors of a $\sigma\pi$ into a type and write a view type that allows us to look at an inhabitant of a sum of products as a constructor and data.

First, we define constructors:

```
\begin{array}{l} \operatorname{cons}\#:\left\{n:\,\mathbb{N}\right\}\to\sigma\pi\ n\to\mathbb{N}\\ \operatorname{cons}\#=\operatorname{length} \end{array} \begin{array}{l} \operatorname{Constr}:\left\{n:\,\mathbb{N}\right\}(ty:\,\sigma\pi\ n)\to\operatorname{Set}\\ \operatorname{Constr}\ ty=\operatorname{Fin}\ (\operatorname{cons}\#\ ty) \end{array}
```

Now, a constructor of type C expects some arguments to be able to make an element of type C. This is a product, we call it the typeOf the constructor.

```
\begin{array}{l} \mathsf{typeOf}: \{n: \mathbb{N}\}(ty: \sigma\pi\ n) \to \mathsf{Constr}\ ty \to \pi\ n \\ \mathsf{typeOf}\ []\ () \\ \mathsf{typeOf}\ (x:: ty)\ \mathsf{fz} = x \\ \mathsf{typeOf}\ (x:: ty)\ (\mathsf{fs}\ c) = \mathsf{typeOf}\ ty\ c \end{array} Injecting is fairly simple.
```

```
\begin{split} & \text{inject} : \{n : \mathbb{N}\} \{A : \mathsf{Parms}\ n\} \{ty : \sigma\pi\ n\} \\ & \to (i : \mathsf{Constr}\ ty) \to [\![\![\ \mathsf{typeOf}\ ty\ i\ ]\!]_p\ A \\ & \to [\![\![\ ty\ ]\!]\!]\ A \\ & \text{inject}\ \{ty = [\![\!]\!]\}\ ()\ cs \\ & \text{inject}\ \{ty = x :: ty\}\ \mathsf{fz}\ cs \ = \mathsf{i1}\ cs \\ & \text{inject}\ \{ty = x :: ty\}\ (\mathsf{fs}\ i)\ cs \ = \mathsf{i2}\ (\mathsf{inject}\ i\ cs) \end{split}
```

We finish off with a *view* of $[ty]_A$ as a constructor and some data. This greatly simplify the algorithms later on.

```
\begin{array}{l} \mathsf{data} \; \mathsf{SOP} \; \{n : \; \mathbb{N}\} \{A : \mathsf{Parms} \; n\} \{ty : \sigma \pi \; n\} \; \colon \llbracket \; ty \; \rrbracket \; A \to \mathsf{Set} \; \mathsf{where} \\ \mathsf{strip} \; \colon (i : \; \mathsf{Constr} \; ty) (\mathit{ls} : \; \llbracket \; \mathsf{typeOf} \; ty \; i \; \rrbracket_p \; A) \\ \to \mathsf{SOP} \; (\mathsf{inject} \; i \; \mathit{ls}) \end{array}
```

1.2 Agda Details

Here we clarify some Agda specific details that are agnostic to the big picture. This can be safely skipped on a first iteration.

As we mentioned above, our codes represent functors on n variables. Obviously, to program with them, we need to apply these to something. The denotation receives a function $\operatorname{\mathsf{Fin}} n \to \operatorname{\mathsf{Set}}$, denoted $\operatorname{\mathsf{Parms}} n$, which can be seen as a valuation for each type variable.

In the following sections, we will be dealing with values of $[ty]_A$ for some class of valuations A, though. We need to have decidable equality for A k and some mapping from A k to \mathbb{N} for all k. We call such valuations a well-behaved parameter:

```
record WBParms \{n: \mathbb{N}\}(A: \mathsf{Parms}\ n): \mathsf{Set}\ \mathsf{where} constructor wb-parms field \mathsf{parm\text{-}size}: \ \forall \{k\} \to A\ k \to \mathbb{N} \mathsf{parm\text{-}cmp}: \ \forall \{k\}(x\ y: A\ k) \to \mathsf{Dec}\ (x \equiv y)
```

TODO

The field parm-size is not really needed anymore! Remove it!

The following sections discuss functionality that does not depend on parameters to codes. Hence, we will be passing them as Agda module parameters. We also set up a number of synonyms to already fix the aforementioned parameter. The first diffing technique we discuss is the trivial diff. It's module is declared as follows:

We stick to this nomenclature throughout the code. The first line handles constant types: ks# is how many constant types we have, ks is the vector of such types and keqs is an indexed vector with a proof of decidable equality over such types. The second line handles type parameters: parms# is how many type-variables our codes will have, A is the valuation we are using and WBA is a proof that A is $well\ behaved$.

TODO

Now parameters are setoids, we can drop out the WBA record.

Below are the synonyms we use for the rest of the code:

```
\llbracket \_ \rrbracket_a : \mathsf{Atom} \to \mathsf{Set}
[ a ]_a = interp_a \ a \ A
 \llbracket \_ \rrbracket : \mathsf{U} \to \mathsf{Set} \\ \llbracket \ u \ \rrbracket = \mathsf{interp}_s \ u \ A 
UUSet: Set<sub>1</sub>
\mathsf{UUSet} = \mathsf{U} \to \mathsf{U} \to \mathsf{Set}
\mathsf{AASet} : \mathsf{Set}_1
\mathsf{AASet} = \mathsf{Atom} \to \mathsf{Atom} \to \mathsf{Set}
\Pi\Pi Set : Set_1
\Pi\Pi\mathsf{Set}=\Pi\to\Pi\to\mathsf{Set}
contr : \forall \{a \ b\} \{A : \mathsf{Set} \ a\} \{B : \mathsf{Set} \ b\}
      \rightarrow (A \rightarrow A \rightarrow B) \rightarrow A \rightarrow B
\mathsf{contr}\; p\; x = p\; x\; x
UU \rightarrow AA : UUSet \rightarrow AASet
UU \rightarrow AA P a a' = P (\alpha a) (\alpha a')
\rightarrow\!\!\alpha: \{a: \mathsf{Atom}\} \rightarrow \llbracket \ a \ \rrbracket_a \rightarrow \llbracket \ \alpha \ a \ \rrbracket
\rightarrow \alpha \ k = i1 \ (k \ , unit)
```

2 Computing and Representing Patches

Intuitively, a Patch is some description of a transformation. Setting the stage, let A and B be a types, x:A and y:B elements of such types. A patch between x and y must specify a few parts:

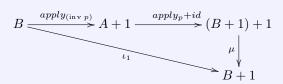
- i) An $apply_p: A \to Maybe\ B$ function,
- ii) such that $apply_p x \equiv \text{just } y$.

Well, $apply_p$ can be seen as a functional relation (R is functional iff $img \ R \subseteq id$) from A to B. We call this the "application" relation of the patch, and we will denote it by $p^{\flat} \subseteq A \times B$.

Needs discussion:

There is still a lot that could be said about this. I feel like p^{\flat} should also be invertible in the sense that:

- i) Let (inv p) denote the inverse patch of p, which is a patch from B to A.
- ii) Then, $p^{\flat} \cdot (\text{inv } p)^{\flat} \subseteq id$ and $(\text{inv } p)^{\flat} \cdot p^{\flat} \subseteq id$, Assuming $(\text{inv } p)^{\flat}$ is also functional, we can use the maybe monad to represent these relations in **Set**. Writing the first equation on a diagram in **Set**, using the apply functions:



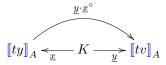
iii) This is hard to play ball with. We want to say, in a way, that $x(p^{\flat})$ y iff $y((\text{inv }p)^{\flat})$ x. That is, (inv p) is the actual inverse of p. Using relations, one could then say that $(\text{inv }p)^{\flat}$ is the converse of (p^{\flat}) . That is: $(\text{inv }p)^{\flat} \equiv (p^{\flat})^{\circ}$. But, if $(\text{inv }p)^{\flat}$ is functional, so is $(p^{\flat})^{\circ}$. This is the same as saying that p^{\flat} is entire! If p^{\flat} is functional and entire, it is a function (and hence, total!). And that is not true.

Now, let us discuss some code and build some intuition for what is what in the above schema. We will present different parts of the code, how do they relate to this relational view and give examples here and there!

2.1 Trivial Diff

The simplest possible way to describe a transformation is to say what is the source and what is the destination of such transformation. This can be accomplished by the Diagonal functor, Δ , just fine.

Now, take an element (x, y): $\Delta ty tv$. The "application" relation it defines is trivial: $\{(x,y)\}$, or, in PF style:



Where, for any $A, B \in Set$ and $x : A, \underline{x} \subseteq A \times B$ represents the everywhere x relation, defined by

$$x = \{(x, b) \mid b \in B\}$$

This is a horrible patch however: We can't calculate with it because we don't know anything about how x changed into y. Note, however, that $(x, y)^{\flat} \equiv \underline{y} \cdot \underline{x}^{\circ}$ is trivially functional.

Needs discussion:

In the code, we actually define the "application" relation of Δ as:

$$(x, x)^{\flat} = id$$
$$(x, y)^{\flat} = y \cdot \underline{x}^{\circ}$$

This suggests that copies might be better off being handled by the trivial diff. We will return to this discussion in section ??

5

2.1.1 Trivial Diff, in Agda

We will be using Δ ty tv for the three levels of our universe: atoms, products and sums. We distinguish between the different Δ 's with subscripts a, p and s respectively. They only differ in type. The treatment they receive in the code is exactly the same! Below is how they are defined:

```
\begin{array}{l} \mathsf{delta} : \forall \{a\} \{A : \mathsf{Set}\ a\} (P : A \to \mathsf{Set}) \\ \to A \to A \to \mathsf{Set} \\ \mathsf{delta}\ P\ a_1\ a_2 = P\ a_1 \times P\ a_2 \end{array}
```

Hence, we define $\Delta_x = \mathsf{delta} \, \llbracket \cdot \rrbracket_{\mathsf{x}}$, for $x \in \{a, p, s\}$.

2.2 Spines

We can make the trivial diff better by identifying whether or not x and y agree on something! In fact, we will aggresively look for copying oportunities. We start by checking if x and y are, in fact, equal. If they are, we say that the patch that transforms x into y is copy. If they are not equal, they might have the same constructor. If they do, the say that the constructor is copied and we put the data side by side (zip). Otherwise, there is nothing we can do on this phase and we just return Δx .

Note that the *spine* forces x and y to be of the same type! In practice, we are only interested in diffing elements of the same language. It does not make sense to diff a C source file against a Haskell source file.

Nevertheless, we define an S structure to capture this longest common prefix of x and y; which, for the SoP universe is very easy to state.

```
\begin{array}{l} \mathsf{data} \; \mathsf{S} \; (P : \mathsf{UUSet}) : \; \mathsf{U} \to \mathsf{Set} \; \mathsf{where} \\ \mathsf{SX} \; : \; \{ty : \; \mathsf{U}\} \to P \; ty \; ty \to \mathsf{S} \; P \; ty \\ \mathsf{Scp} \; : \; \{ty : \; \mathsf{U}\} \to \mathsf{S} \; P \; ty \\ \mathsf{Scns} : \; \{ty : \; \mathsf{U}\}(i : \; \mathsf{Constr} \; ty) \\ \to \mathsf{ListI} \; (\mathsf{contr} \; P \circ \alpha) \; (\mathsf{typeOf} \; ty \; i) \\ \to \mathsf{S} \; P \; ty \end{array}
```

Remember that contr P x = P x x and $\alpha : \text{Atom } n \to \sigma \pi n$; Here, Listl P l is an indexed list where the elements have type $P l_i$, for every $l_i \in l$. We will treat this type like an ordinary list for the remainder of this document.

Note that S makes a functor (actually, a free monad!) on P, and hence, we can map over it:

```
\begin{array}{lll} \mathsf{S-map} &: \{ty : \mathsf{U}\} \\ &\quad \{P\ Q : \mathsf{UUSet}\}(X : \forall \{k\ v\} \rightarrow P\ k\ v \rightarrow Q\ k\ v) \\ &\quad \rightarrow \mathsf{S}\ P\ ty \rightarrow \mathsf{S}\ Q\ ty \\ \mathsf{S-map}\ f\ (\mathsf{SX}\ x) &= \mathsf{SX}\ (f\ x) \\ \mathsf{S-map}\ f\ \mathsf{Scp} &= \mathsf{Scp} \\ \mathsf{S-map}\ f\ (\mathsf{Scns}\ i\ xs) &= \mathsf{Scns}\ i\ (\mathsf{map}_i\ f\ xs) \end{array}
```

Computing a spine is easy, first we check whether or not x and y are equal. If they are, we are done. If not, we look at x and y as true sums of products and check if their constructors are equal, if they are, we zip the data together. If they are not, we zip x and y together and give up.

```
\begin{array}{lll} \operatorname{spine-cns}: \{ty: \ \mathbb{U}\}(x\ y: \ \mathbb{I}\ ty\ \mathbb{J}) \to \operatorname{S}\ \Delta_s\ ty \\ \operatorname{spine-cns}\ x\ y & \text{with sop}\ x\ | \operatorname{sop}\ y \\ \operatorname{spine-cns}\ _-\ | & \operatorname{strip}\ cx\ dx\ | & \operatorname{strip}\ cy\ dy \\ & \text{with}\ cx\stackrel{?}{=}-\operatorname{Fin}\ cy \\ \ldots | & \operatorname{no}\ _-\  = \operatorname{SX}\ (\operatorname{inject}\ cx\ dx\ , \operatorname{inject}\ cy\ dy) \\ \operatorname{spine-cns}\ _-\ | & \operatorname{strip}\ _-\ dx\ | & \operatorname{strip}\ cy\ dy \\ | & \operatorname{yes}\ \operatorname{refl}\  = \operatorname{Scns}\ cy\ (\operatorname{zip}_p\ dx\ dy) \\ & \operatorname{spine}: \{ty: \ \mathbb{U}\}(x\ y: \ \mathbb{I}\ ty\ \mathbb{J}) \to \operatorname{S}\ \Delta_s\ ty \\ & \operatorname{spine}\ \{ty\}\ x\ y \\ & \text{with}\ \operatorname{dec-eq}\ _-\ _-\ ^?-\operatorname{A}\ _-\ ty\ x\ y \\ \ldots | & \operatorname{yes}\ _-\  = \operatorname{Scp}\ \ldots | & \operatorname{no}\ _-\  = \operatorname{spine-cns}\ x\ y \\ & \operatorname{zip}_p: \{ty: \ \Pi\} \\ & \to \ \|\ ty\ \|_p \to \ \|\ ty\ \|_p \to \operatorname{Listl}\ (\lambda\ k\to \Delta_s\ (\alpha\ k)\ (\alpha\ k))\ ty \\ & \operatorname{zip}_p\ \{\mathbb{J}\}\ _-\ _-\  = \mathbb{J}\ \\ & \operatorname{zip}_p\ \{_-\ ::\ ty\}\ (x\ ,\ xs)\ (y\ ,\ ys) \\ & = \ (\operatorname{il}\ (x\ ,\ \operatorname{unit})\ ,\ \operatorname{il}\ (y\ ,\ \operatorname{unit}))\ ::\ \operatorname{zip}_p\ xs\ ys \\ \end{array}
```

The "application" relations specified by a spine $s=\operatorname{spine} x\ y,$ denoted s^{\flat} are defined by:

where inj_i is the injection, with constructor i, into $\coprod_k T_k$. It corresponds to the relational lifting of function injection.

Note that, in the (SX p) case, we simply ask for the "application" relation of p. The algorithm produces a S Δ_s , so we have pairs on the leaves of the spine. In fact, either we have only one leave or we have $arity\ C_i$ leaves, where C_i is the common constructor of x and y in spine x y.

For a running example, let's consider a datatype defined by:

```
2-3-TREE-F : \sigma\pi 1
2-3-TREE-F = []
\oplus (K kN) \otimes I \otimes I \otimes []
\oplus (K kN) \otimes I \otimes I \otimes I \otimes []
\oplus []
```

We ommit the fz for the I parts, as we only have one type variable. We also use $_\oplus_$ and $_\otimes_$ as aliases for $_::_$ with different precedences. As expected, there are three constructors:

```
2-node' 3-node' nil' : Constr 2-3-TREE-F
nil' = fz
2-node' = fs fz
3-node' = fs (fs fz)
```

We can then consider a few spines over $[2-3-TREE-F]_{Unit}$ to illustrate the algorithm:

In the case where the spine is Scp or Scns i there is nothing left to be done and we have the best possible diff. Note that on the Scns i case we do not allow for rearanging of the parameters of the constructor i.

In the case where the spine is SX, we can do a better job! We can record which constructor changed into which and try to reconcile the data from both the best we can. Going one step at a time, let's first change one constructor into the other.

2.3 Constructor Changes

Let's take an example where the spine can not copy anything:

```
s = \text{spine} (2\text{-node}' \ 10 \ \text{unit unit}) (3\text{-node}' \ 10 \ \text{unit unit unit})
= SX (2-node' \( 10 \ \text{unit unit unit} \), \( 3\text{-node}' \ 10 \ \text{unit unit unit} \)
```

Here, we wish to say that we changed a 2-node into a 3-node. But we are then left with a problem about what to do with the data inside the 2-node and 3-node; this is where the notion of alignment will be in the picture. For now, we abstract it away by the means of a parameter, just like we did with the S. This time, however, we need something that receives products as inputs.

```
\begin{array}{l} \mathsf{data} \ \mathsf{C} \ (P: \Pi\Pi \mathsf{Set}) : \mathsf{U} \to \mathsf{U} \to \mathsf{Set} \ \mathsf{where} \\ \mathsf{CX} \ : \{ty \ tv : \ \mathsf{U}\} \\ \ \to (i: \mathsf{Constr} \ ty)(j: \mathsf{Constr} \ tv) \\ \ \to P \ (\mathsf{typeOf} \ ty \ i) \ (\mathsf{typeOf} \ tv \ j) \\ \ \to \mathsf{C} \ P \ ty \ tv \end{array}
```

Note that C also makes up a functor, and hence can be mapped over:

```
\begin{array}{ll} \mathsf{C\text{-}map} & : & \{ty \; tv : \; \mathsf{U}\} \\ & & \{P \; Q : \mathsf{\Pi} \mathsf{\Pi} \mathsf{Set}\}(X : \forall \{k \; v\} \; \rightarrow \; P \; k \; v \; \rightarrow \; Q \; k \; v) \\ & & \rightarrow \; \mathsf{C} \; P \; ty \; tv \; \rightarrow \; \mathsf{C} \; Q \; ty \; tv \\ \mathsf{C\text{-}map} \; f \; (\mathsf{CX} \; i \; j \; x) \; = \; \mathsf{CX} \; i \; j \; (f \; x) \end{array}
```

Computing an inhabitant of C is trivial:

```
\begin{array}{ll} \mathsf{change} : \{ ty \ tv : \ \mathsf{U} \} \to \llbracket \ ty \ \rrbracket \to \llbracket \ tv \ \rrbracket \to \mathsf{C} \ \Delta_p \ ty \ tv \\ \mathsf{change} \ x \ y \ \mathsf{with} \ \mathsf{sop} \ x \mid \mathsf{sop} \ y \\ \mathsf{change} \ \mid \ \mathsf{strip} \ cx \ dx \mid \mathsf{strip} \ cy \ dy = \mathsf{CX} \ cx \ cy \ (dx \ , \ dy) \end{array}
```

Now that we can compute change of constructors, we can refine our s above. We can compute $\mathsf{S-map}$ change s and we will have:

```
S-map change s = SX (CX 2\text{-node}' 3\text{-node}' ((10, unit, unit), (10, unit, unit, unit)))
```

The "application" relation induced by C is trivial:

```
V \xleftarrow{(\mathsf{CX}\ i\ j\ p)^{\flat}} T \\ \inf_{\mathsf{j}_j} \bigwedge^{\flat} \bigvee_{\mathsf{tinj}_i^{\circ}} \mathsf{typeOf}\ V\ j \xleftarrow{p^{\flat}} \mathsf{typeOf}\ T\ i
```

Note that up until now, everything was deterministic! This is something we are bound to lose when talking about alignment.

2.4 Aligning Everything

Following a similar reasoning as from S to C; the leaves of a C produced through change will NEVER contain a coproduct as the topmost type. Hence, we know that they will contain either a product, or a constant type, or a type variable. In the case of a constant

type or a type variable, there is not much we can do at the moment, but for a product we can refine this a little bit more before using Δ^1 .

¹In fact, splitting the different stages of the algorithm into different types reinforced our intuition that the alignment is the source of difficulties. As we shall see, we now need to introduce non-determinism.