

Best Title in the Universe

42

Victor Cacciari Miraldo Wouter Swierstra

University of Utrecht

{v.cacciarimiraldo,w.s.swierstra} at uu.nl

Abstract

stuff

Categories and Subject Descriptors D.1.1 [look]: for—this

General Terms Haskell

Keywords Haskell

1. Introduction

The majority of version control systems handle patches in a non-structured way. They see a file as a list of lines that can be inserted, deleted or modified, with no regard to the semantics of that specific file. The immediate consequence of such design decision is that we, humans, have to solve a large number of conflicts that arise from, in fact, non conflicting edits. Implementing a tool that knows the semantics of any file we happen to need, however, is no simple task, specially given the plethora of file formats we see nowadays.

This can be seen from a simple example. Lets imagine Alice and Bob are iterating over a cake’s recipe. They decide to use a version control system and an online repository to keep track of their modifications.

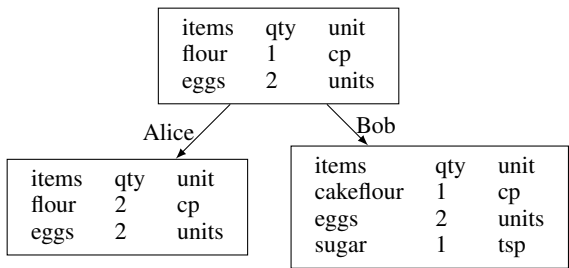


Figure 1. Sample CSV files

Lets say that both Bob and Alice are happy with their independent changes and want to make a final recipe. The standard way to track differences between files is the `diff3` [FIXBIB] unis tool. Running `diff3 Alice.csv 0.csv Bob.csv` would result in the output presented in figure 2. Every tag `====` marks a difference. Three

locations follows, formatted as `file:line` type. The change type can be a *Change*, *Append* or *Delete*. The first one, says that file 1 (`Alice.csv`) has a change in line 2 (`1:2c`) which is `flour, 2 , cp`; and files 2 and 3 have different changes in the same line. The tag `====3` indicates that there is a difference in file 3 only. Files 1 and 2 should append what changed in file 3 (line 4).

```
====
1:2c
    flour, 2 , cp
2:2c
    flour, 1 , cp
3:2c
    cakeflour, 1 , cp
====3
1:3a
2:3a
3:4c
    sugar, 1 , tsp
```

Figure 2. Output from `diff3`

If we try to merge the changes, `diff3` will flag a conflict and therefore require human interaction to solve it, as we can see by the presence of the `====` indicator in its output. However, Alice’s and Bob’s edits, in figure 1 do *not* conflict, if we take into account the semantics of CSV files. Although there is an overlapping edit at line 1, the fundamental editing unit is the cell, not the line.

We propose a structural diff that is not only generic but also able to track changes in a way that the user has the freedom to decide which is the fundamental editing unit. Our work was inspired by [?] and [?]. We did extensive changes in order to handle structural merging of patches. We also propose extensions to this algorithm capable of detecting purely structural operations such as refactorings and cloning.

The paper begins by exploring the problem, generically, in the Agda [FIXBIB] language. Once we have a provably correct algorithm, the details of a Haskell implementation of generic diff’ing are sketched. To open ground for future work, we present a few extensions to our initial algorithm that could be able to detect semantical operations such as *cloning* and *swapping*.

Contributions

- *Study of a more algebraic patch theory.*
- *Agda model.*
- *Haskell Prototype.*

Background

Permission to make digital or hard copies of part or all 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. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, or to redistribute to lists, contact the Owner/Author. Request permissions from permissions@acm.org or Publications Dept., ACM, Inc., fax +1 (212) 869-0481. Copyright held by Owner/Author. Publication Rights Licensed to ACM.

- Should we have this section? It could be nice to at least mention the edit distance problem and that in the untyped scenario, the best running time is of $O(n^3)$. Types should allow us to bring this time lower.

2. Structural Diffing

Alice and Bob were both editing a CSV file which represents data that is isomorphic to $[[Atom\ String]]$, where *Atom* *a* is a simple tag that indicates that *as* should be treated as atomic.

- What do we mean by structural?
- Give some context: Tree-edit distance;
- We seek to obtain a system with something close to residuals.

To Research!

- Check out the antidiagonal with more attention: <http://blog.sigfpe.com/2007/09/type-of-distinct-pairs.html>
- The LCS problem is closely related to diffing. We want to preserve the LCS of two structures! How does our diffing relate? Does this imply maximum sharing?

2.1 Context Free Datatypes

- Explain the universe we're using.
- Explain the intuition behind our *D* datatype.
- Mention that it is correct.

Took from [?].

```
data U : N → Set where
  u0  : {n : N} → U n
  u1  : {n : N} → U n
  ⊕_  : {n : N} → U n → U n → U n
  ⊗_  : {n : N} → U n → U n → U n
  β_  : {n : N} → U (suc n) → U n → U n
  μ_  : {n : N} → U (suc n) → U n
  vl_ : {n : N} → U (suc n)
  wk_ : {n : N} → U n → U (suc n)
```

- Explain the patching problem.
- We want a type-safe approach.
- Argue that the types resulting from our parser are in a sub-language of what we treated next.
- introduce edit-script, diffing and patching or apply

2.2 Patches over a Context Free Type

- Explain that a patch is something which we can apply.
- Loh's approach is too generic, as the diff function should have type $a \rightarrow a \rightarrow D\ a$.

In order to simplify the presentation, we are gonna explicitly name variables and write our types in a more mathematical fashion, other than the Agda encoding. As we discussed earlier, a patch is an object that track differences in a given type. Different types will allow for different types of changes.

Definition 2.1 (Simple Patch). We define a (simple) patch *D ty* by induction on *ty* as:

$$\begin{aligned}
 D\ 0 &= 0 \\
 D\ 1 &= 1 \\
 D\ (x \times y) &= D\ x \times D\ y \\
 D\ (x + y) &= (D\ x + D\ y) + 2 \times (x \times y) \\
 D\ (\mu X.F\ X) &= \mu X.(1 \\
 &\quad + D\ (F\ 1) \times X \\
 &\quad + 2 \times (F\ 1) \times X \\
 &\quad)
 \end{aligned}$$

Where 1 and 0 are the usual terminal and initial objects of a given category.

Let's see the coproduct case in more detail. There are four different possibilities for the changes seen in a coproduct, just like there are four different combinations of constructors for two objects of type *Either a b*. The first and second options, namely *D x* and *D y* track differences of a *Left a* into a *Left a'* and a *Right b* into a *Right b'*, respectively. The other possibilities are representing a *Left a* becoming a *Right b* or vice-versa. The other branches are straight-forward.

Producing Patches Definition 2.1 provides us with some intuition on how one would define patches for a given datatype. The actual definition (figure 4) is more complicated, though. The *Diff* type takes one parameter, used to give a free-monad [FIXBIB] structure, and two indexes which indicate the type for which that diff is intended. We then define:

```
Patch : {n : N} → Tel n → U n → Set
Patch t ty = D ⊥p t ty
```

Where $\perp_p = \lambda_ \rightarrow \perp$.

Our first goal is to produce patches, or to differentiate between to objects of the same type. We can do that generically through

```
gdiff : {n : N} {t : Tel n} {ty : U n}
       → EIU ty t → EIU ty t → Patch t ty
gdiff {ty = vl} (top a) (top b) = D-top (gdiff a b)

gdiff {ty = wk u} (pop a) (pop b) = D-pop (gdiff a b)

gdiff {ty = β F x} (red a) (red b) = D-β (gdiff a b)

gdiff {ty = u1} void void = D-void

gdiff {ty = ty ⊗ tv} (ay , av) (by , bv)
  = D-pair (gdiff ay by) (gdiff av bv)

gdiff {ty = ty ⊕ tv} (inl ay) (inl by) = D-inl (gdiff ay by)
gdiff {ty = ty ⊕ tv} (inr av) (inr bv) = D-inr (gdiff av bv)
gdiff {ty = ty ⊕ tv} (inl ay) (inr bv) = D-setl ay bv
gdiff {ty = ty ⊕ tv} (inr av) (inl by) = D-setr av by

gdiff {ty = μ ty} a b = D-mu (gdiffL (a :: []) (b :: []))
```

- wrap it up?

Definition 2.2 (Defined). We say that a patch p_a is defined for an input *a* iff there exists an object a' such that:

$$\text{apply } p_a\ a \equiv \text{Just } a'$$

Fixed Points The treatment for fixed points has to be made uniform, somehow, if we want a generic algorithm by the end of the day. What makes fixed points different than regular algebraic types is that they can grow or shrink arbitrarily, and our diff function has to take that into account.

Recalling the fixed point clause of simple patches (def 2.1),

$$\begin{aligned} D(\mu X.F X) &= \mu X.(1 \\ &+ D(F 1) \times X \\ &+ 2 \times (F 1) \times X \\ &) \end{aligned}$$

it is straight forward to see that the $D(\mu X.F X)$ is isomorphic to a list with three recursive constructors and a non-recursive one. Following the edit operations studied by Löh[?], we have an *insert*, a *delete* and a *end* edit operations. The big difference is that instead of copying, we have a constructor that track changes inside a constructor of $\mu X.F X$, we call this a *down* edit operation.

We heavily rely on the fact that $\mu X.F X \approx F 1 \times [\mu X.F X]$, that is, any inhabitant of a fixed-point type can be seen as a non-recursive head and a list of recursive children, or, expressed in our generic setting:

$\text{Open}\mu : \{n : \mathbb{N}\} \rightarrow \text{Tel } n \rightarrow \text{U}(\text{succ } n) \rightarrow \text{Set}$
 $\text{Open}\mu \text{ } t \text{ } ty = \text{EIU } ty (\text{tcons u1 } t) \times \text{List}(\text{EIU}(\mu \text{ } ty) \text{ } t)$

$\mu\text{-open} : \{n : \mathbb{N}\} \{t : \text{Tel } n\} \{ty : \text{U}(\text{succ } n)\}$
 $\rightarrow \text{EIU}(\mu \text{ } ty) \text{ } t \rightarrow \text{Open}\mu \text{ } t \text{ } ty$

$\mu\text{-close} : \{n : \mathbb{N}\} \{t : \text{Tel } n\} \{ty : \text{U}(\text{succ } n)\}$
 $\rightarrow \text{Open}\mu \text{ } t \text{ } ty \rightarrow \text{Maybe}(\text{EIU}(\mu \text{ } ty) \text{ } t \times \text{List}(\text{EIU}(\mu \text{ } ty) \text{ } t))$

Although we could have used vectors of a fixed length and made this a total isomorphism, we would have more problems than benefits. This will be discussed in section 2.5. Nonetheless, an important soundness result has been proven:

$\mu\text{-close-resp-arity}$
 $: \{n : \mathbb{N}\} \{t : \text{Tel } n\} \{ty : \text{U}(\text{succ } n)\} \{a : \text{EIU}(\mu \text{ } ty) \text{ } t\}$
 $\{hdA : \text{EIU } ty (\text{tcons u1 } t)\} \{chA l : \text{List}(\text{EIU}(\mu \text{ } ty) \text{ } t)\}$
 $\rightarrow \mu\text{-open } a \equiv (hdA, chA)$
 $\rightarrow \mu\text{-close}(hdA, chA ++ l) \equiv \text{just}(a, l)$

2.3 The Cost Function

2.4 Sharing of Recursive Subterms

- If we want to be able to share recursive subexpressions we need a mutually recursive approach.
- Or, this will be handled during conflict solving. See refactoring.

2.5 Remarks on Type Safety

- only the interface to the user can be type-safe, otherwise we don't have our free-monad multiplication.

3. Patch Propagation

Let's say Bob and Alice perform edits in a given object, which are captured by patches p and q , shown in figure 3. The natural question to ask is *how do we join these changes*.

The residual p/q of two patches p and q only makes sense if both p and q are aligned, that is, are defined for the same input. It captures the notion of incorporating the changes made by p in an object that has already been modified by q .

- *Pijul has this notion of handling a merge as a pushout, but it uses the free co-completion of a rather simple category. This doesn't give enough information for structured conflict solving.*

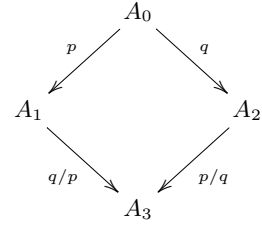


Figure 3. Residual Square

• BACK THIS UP!

In an ideal world, we would expect the residual function to have type $D a \rightarrow D a \rightarrow \text{Maybe}(D a)$, where the partiality comes from receiving two non-aligned patches.

But what if Bob and Alice changes the same cell in their CSV file? Then it is obvious that someone (human) have to chose which value to use in the final, merged, version.

Here we see that this residual operation is where conflicts are introduced in our theory.

3.1 Incorporating Conflicts

In order to track down these conflicts we need a more expressive patch data structure. We chose to parametrize D over an abstract type and add another dummy constructor to it, much like we would do in a free monad construction. The actual data structure we use is presented in figure 4.

- *Show where conflicts arise and the two types of conflicts we identify.*

Note that the first constructor of D just asks for a suitably indexed A . With this in mind, we can now start to define our residual operation.

- *Define Patch = D ⊥*

$_ / _ : \{n : \mathbb{N}\} \{t : \text{Tel } n\} \{ty : \text{U } n\}$
 $\rightarrow \text{Patch } t \text{ } ty \rightarrow \text{Patch } t \text{ } ty \rightarrow \text{Maybe}(D C t \text{ } ty)$

It is interesting to note that this residual operation is somewhat symmetric:

residual-symmetry-thm
 $: \{n : \mathbb{N}\} \{t : \text{Tel } n\} \{ty : \text{U } n\} \{k : D C t \text{ } ty\}$
 $\rightarrow (d1 \text{ } d2 : \text{Patch } t \text{ } ty)$
 $\rightarrow d1 / d2 \equiv \text{just } k$
 $\rightarrow \Sigma(D C t \text{ } ty \rightarrow D C t \text{ } ty)$
 $(\lambda op \rightarrow d2 / d1 \equiv \text{just}(D\text{-map } C\text{-sym}(op \text{ } k)))$

residual-sym-stable : $\{n : \mathbb{N}\} \{t : \text{Tel } n\} \{ty : \text{U } n\} \{k : D C t \text{ } ty\}$
 $\rightarrow (d1 \text{ } d2 : \text{Patch } t \text{ } ty)$
 $\rightarrow d1 / d2 \equiv \text{just } k$
 $\rightarrow \text{forget } \langle M \rangle (d2 / d1) \equiv \text{just}(\text{map}(\downarrow\text{-map-}\downarrow C\text{-sym})(\text{forget } k))$

This means that p/q and q/p , although different, have the same conflicts (up to symmetry).

3.2 Solving Conflicts

- *This is highly dependent on the structure.*
 - *some structures might allow permutations, refactorings, etc... whereas others might not.*
- *How do we go generic? Free-monads to the rescue!*

```

mutual
data D {a} (A : {n : N} → Tel n → U n → Set a) : {n : N} → Tel n → U n → Set a where
  D-A : {n : N} {t : Tel n} {ty : U n} → A t ty → D A t ty

  D-void : {n : N} {t : Tel n} → D A t u1
  D-inl : {n : N} {t : Tel n} {a b : U n}
    → D A t a → D A t (a ⊕ b)
  D-inr : {n : N} {t : Tel n} {a b : U n}
    → D A t b → D A t (a ⊕ b)
  D-setl : {n : N} {t : Tel n} {a b : U n}
    → EU a t → EU b t → D A t (a ⊕ b)
  D-setr : {n : N} {t : Tel n} {a b : U n}
    → EU b t → EU a t → D A t (a ⊕ b)
  D-pair : {n : N} {t : Tel n} {a b : U n}
    → D A t a → D A t b → D A t (a ⊗ b)
  D-mu : {n : N} {t : Tel n} {a : U (suc n)}
    → List (Dμ A t a) → D A t (μ a)
  D-β : {n : N} {t : Tel n} {F : U (suc n)} {x : U n}
    → D A (tcons x t) F → D A t (β F x)
  D-top : {n : N} {t : Tel n} {a : U n}
    → D A t a → D A (tcons a t) v1
  D-pop : {n : N} {t : Tel n} {a b : U n}
    → D A t b → D A (tcons a t) (wk b)

data Dμ {a} (A : {n : N} → Tel n → U n → Set a) : {n : N} → Tel n → U (suc n) → Set a where
  Dμ-A : {n : N} {t : Tel n} {a : U (suc n)} → A t (μ a) → Dμ A t a
  Dμ-ins : {n : N} {t : Tel n} {a : U (suc n)} → ValU a t → Dμ A t a
  Dμ-del : {n : N} {t : Tel n} {a : U (suc n)} → ValU a t → Dμ A t a
  Dμ-cpy : {n : N} {t : Tel n} {a : U (suc n)} → ValU a t → Dμ A t a
  Dμ-dwn : {n : N} {t : Tel n} {a : U (suc n)} → ValU a t → D A t (β a u1) → Dμ A t a

```

Figure 4. Complete Definition of D

4. Sketching a Control Version System

- Different views over the same datatype will give different diffs.
- `newtype` annotations can provide a great bunch of control over the algorithm.
- Directories are just rosetrees...

5. Related Work

- People have done similar things... or not.

6. Conclusion

- This is what we take out of it.