

# Type-directed Diffing of Structured Data

Victor Cacciari Miraldo, Pierre-Évariste Dagand  
and Wouter Swierstra

January 5, 2018



# The diff utility

The Unix `diff` utility compares two files line-by-line, computing the smallest number of insertions and deletions to transform one into the other.

It was developed as far back as 1976 – but still forms the heart of many modern version control systems such as `git`, `mercurial`, `svn`, and many others.



## Example: comparing two files

jabber.txt

```
Twas brillig, and the slithy toves  
Waved to Mars, where a robot roves;  
Did gyre and gimble in the wabe;  
And the mome raths outgrabe.
```

wocky.txt

```
Twas brillig, and the slithy toves  
Did gyre and gimble in the wabe;  
All mimsy were the borogoves,  
And the mome raths outgrabe.
```



## Example: comparing two files

```
Twas brillig, and the slithy toves  
- Waved to Mars, where a robot roves;  
  Did gyre and gimble in the wabe;  
+ All mimsy were the borogoves,  
  And the mome raths outgrabe.
```

The diff utility computes a **patch**, that can be used to transform the one file into the other.



# Smallest edit script

Crucially, `diff` always computes the **smallest** patch – minimizing the number of insertions and deletions.



# Smallest edit script

Crucially, `diff` always computes the **smallest** patch – minimizing the number of insertions and deletions.

Tries to preserve as much information as possible.



# Smallest edit script

Crucially, `diff` always computes the **smallest** patch – minimizing the number of insertions and deletions.

Tries to preserve as much information as possible.

But sometimes it still doesn't do a very good job.



## Example: comma separated values

bibliography.csv

Lewis Carroll, The alphabet cipher

Lewis Carroll, The game of logic

Lewis Carroll, The hunting of the snark

How would this file change if I add publication dates?





## Example: comma separated values

- Lewis Carroll, The alphabet cipher
- + Lewis Carroll, The alphabet cipher, 1868
- Lewis Carroll, The game of logic
- + Lewis Carroll, The game of logic, 1887
- Lewis Carroll, The hunting of the snark
- + Lewis Carroll, The hunting of the snark, 1876



## Example: comma separated values

- Lewis Carroll, The alphabet cipher
- + Lewis Carroll, The alphabet cipher, 1868
- Lewis Carroll, The game of logic
- + Lewis Carroll, The game of logic, 1887
- Lewis Carroll, The hunting of the snark
- + Lewis Carroll, The hunting of the snark, 1876

Syntactically changes **every line**.

Semantically, **data** was not modified.



## Example: comma separated values

- Lewis Carroll, The alphabet cipher
- + Lewis Carroll, The alphabet cipher, 1868
- Lewis Carroll, The game of logic
- + Lewis Carroll, The game of logic, 1887
- Lewis Carroll, The hunting of the snark
- + Lewis Carroll, The hunting of the snark, 1876

Syntactically changes **every line**.

Semantically, **data** was not modified.

Particularly important when diff'ing **source code**.



# What is the diff over structured data?



# Questions

- ▶ How can we represent data types?
- ▶ How can we represent patches on these data types?
- ▶ How can we compute such patches?
- ▶ How can we merge such patches?



# Questions

- ▶ **How can we represent data types?**
- ▶ How can we represent patches on these data types?
- ▶ How can we compute such patches?
- ▶ How can we merge such patches?



# Universe of discourse

We will use Agda as our metalanguage to answer these questions and start by fixing a ‘sums of products’ universe:

```
data Atom : Set where
  K : U → Atom
  I : Atom

Prod : Set
Prod = List Atom

Sum : Set
Sum = List Prod
```

Here we assume some ‘base universe’  $U$ , storing the atomic types such as integers, characters, etc.



# Semantics

We can interpret these types as **pattern functors**:

$$[[\cdot]]_a : \mathbf{Atom} \rightarrow (\mathbf{Set} \rightarrow \mathbf{Set})$$

$$[[I]]_a \quad X = X$$

$$[[K \kappa]]_a \quad X = [[\kappa]]_k$$

$$[[\cdot]]_p : \mathbf{Prod} \rightarrow (\mathbf{Set} \rightarrow \mathbf{Set})$$

$$[[[]]]_p \quad X = \mathbf{Unit}$$

$$[[a :: as]]_p \quad X = [[\alpha]]_a \, X \times [[\pi]]_p \, X$$

$$[[\cdot]]_s : \mathbf{Sum} \rightarrow (\mathbf{Set} \rightarrow \mathbf{Set})$$

$$[[[]]]_s \quad X = \perp$$

$$[[p :: ps]]_s \quad X = [[p]]_p \, X \uplus [[ps]]_s \, X$$





# Fixpoints

Given any element of our ‘sums of products’ universe, we can compute the corresponding pattern functor.

Taking the least fixpoint of this functor allows us to tie the recursive knot:

```
data Fix (s : Sum) : Set where
  ⟨·⟩ :  $\llbracket s \rrbracket_s$  (Fix s) → Fix s
```



## Example: 2-3 trees

We can represent a 2-3-tree, usually defined as follows:

```
data Tree : Set where
  leaf      : Tree
  2-node    : ℕ → Tree → Tree → Tree
  3-node    : ℕ → Tree → Tree → Tree → Tree
```

by the following sum-of-products:

```
TreeF : Sum
TreeF = let leafT      = []
          node2T      = [ K ℕ , | , | ]
          node3T      = [ K ℕ , | , | , | ]
        in [leafT , node2T , node3T ]
```



# Questions

- ▶ How can we represent data types?
- ▶ **How can we represent patches on these data types?**
- ▶ How can we compute such patches?
- ▶ How can we merge such patches?



## 2-3-trees

treeA = 2-node  $\gamma$   $t_1$   $t_2$

treeB = 3-node  $12$  (2-node  $\gamma$   $t_1$  leaf) leaf leaf

What edit script transforms treeA into treeB?



## 2-3-trees

treeA = 2-node  $\gamma$   $t_1$   $t_2$

treeB = 3-node  $12$  (2-node  $\gamma$   $t_1$  leaf) leaf leaf

What edit script transforms treeA into treeB?

It is not just a list of insertions and deletions!

We can insert new constructors, modify values stored in the tree, delete subtrees, or copy over existing data.



# Representing diffs

We define a **type indexed data type**, to account for changes, defining what it means to modify each layer of our universe.

- ▶ sums
- ▶ products
- ▶ atomic values

From these pieces we define our overall type for diffs.



# Spines: changes to sums

Given two arbitrary tree structures,  $x$  and  $y$ , either:

1.  $x$  and  $y$  are equal;
2.  $x$  and  $y$  the same outermost constructor, but are not equal trees;
3.  $x$  and  $y$  have a different outermost constructor.



# Spines: changes to sums

Given two arbitrary tree structures,  $x$  and  $y$ , either:

1.  $x$  and  $y$  are equal;
2.  $x$  and  $y$  the same outermost constructor, but are not equal trees;
3.  $x$  and  $y$  have a different outermost constructor.

Spines,  $S$ , capture these three cases.





Assuming we know what patches on atoms (**At**) and products (**Al**) are, we define:

```
data S ( $\sigma$  : Sum) : Set where
  Scp  : S  $\sigma$ 
  Scns : (C : Constr  $\sigma$ )
        → All At (fields C)
        → S  $\sigma$ 
  Schg : (C1 C2 : Constr  $\sigma$ )
        → Al (fields C1) (fields C2)
        → S  $\sigma$ 
```

Next we define the diff for **products** and **atoms**.



# Alignments: changes to products

How to compare fields of reconciled constructors? (Schg case)



# Alignments: changes to products

How to compare fields of reconciled constructors? (Schg case)

Each value has a **list of fields** – the product structure.

These fields can have different types!



# Alignments: changes to products

How to compare fields of reconciled constructors? (Schg case)

Each value has a **list of fields** – the product structure.

These fields can have different types!

Good news: Unix `diff` algorithm computes a diff for **lists** of lines.



# Alignments: changes to products

To describe a change from one list of constructor fields to another, we require an **edit script** that:

- ▶ changes one field into another;
- ▶ deletes fields;
- ▶ inserts new fields.



# Alignments

data **Al** : Prod → Prod → Set where

**A0** : **Al** [] []

**AX** : **At**  $\alpha$  → **Al**  $\pi_1$   $\pi_2$  → **Al** ( $\alpha :: \pi_1$ ) ( $\alpha :: \pi_2$ )

**Adel** :  $\llbracket \alpha \rrbracket_a$  → **Al**  $\pi_1$   $\pi_2$  → **Al** ( $\alpha :: \pi_1$ )  $\pi_2$

**Ains** :  $\llbracket \alpha \rrbracket_a$  → **Al**  $\pi_1$   $\pi_2$  → **Al**  $\pi_1$  ( $\alpha :: \pi_2$ )

A value of type **Al**  $\pi_1$   $\pi_2$  indicates which fields of one constructor are matched with which fields of another.

Analogous to UNIX `diff` and `lines`.



# Atoms

Finally, we still need to handle our atomic values.

For constant types, we can check if they are equal or not.



# Atoms

Finally, we still need to handle our atomic values.

For constant types, we can check if they are equal or not.

But what about recursive subtrees?





# Handling recursive data types

So far our **spines** compare the outermost constructors.

Oftentimes, one wants to delete certain constructors (exposing its subtrees) or insert new constructors.

We cannot handle such changes with the data types we have seen so far. . .



# Accounting for recursion

Our final patch type identifies three cases:

1. Insertion of a constructor, with a zipper over its fields;
2. Deletion the outermost constructor, with a zipper over its fields;
3. A choice of spine, alignment, and a patch on atomic values;

The first two carry that zipper to point out **where** to insert/delete a subtree. We call this the **context**.



# Applying patches

We can define generic operations – such as patch application – that applies a patch to a given tree:

$$\text{apply} : \text{Patch} \rightarrow \text{Fix } \sigma \rightarrow \text{Maybe } (\text{Fix } \sigma)$$

This patch is guaranteed to **preserve types**.

It may still fail – when encountering an unexpected constructor or atomic value – but it will never produce ill-formed data.



# Questions

- ▶ How can we represent data types?
- ▶ How can we represent patches on these data types?
- ▶ **How can we compute such patches?**
- ▶ How can we merge such patches?



# Computing Patches

Many patches transform one  $x$  into a  $y$ .



# Computing Patches

Many patches transform one  $x$  into a  $y$ .

Too expensive to enumerate. Counting copies not enough.



# Computing Patches

Many patches transform one  $x$  into a  $y$ .

Too expensive to enumerate. Counting copies not enough.

Heuristics to prune the search space.

- ▶ UNIX diff3.
- ▶ Edit scripts preorder traversal.
- ▶ ...



# Computing Patches

Many patches transform one  $x$  into a  $y$ .

Too expensive to enumerate. Counting copies not enough.

Heuristics to prune the search space.

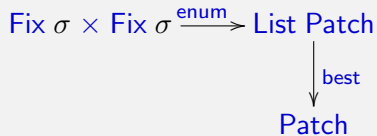
- ▶ UNIX diff3.
- ▶ Edit scripts preorder traversal.
- ▶ ...

Implemented by the means of Oracles.





# Computing Patches: Oracles



# Computing Patches: Oracles

$$\begin{array}{ccc} \text{Fix } \sigma \times \text{Fix } \sigma & \xrightarrow{\text{enum}} & \text{List Patch} \\ \mathcal{O} \downarrow & & \downarrow \text{best} \\ \text{Fix}_a \sigma \times \text{Fix}_a \sigma & \xrightarrow{\text{tr}} & \text{Patch} \end{array}$$

Flag indicating copy:

```
data Fixa (s : Sum) : Set where
  ⟨·,·⟩ : Bool → [[s]]s (Fixa s) → Fixa s
```



# Computing Patches: Oracles

$$\begin{array}{ccc} \text{Fix } \sigma \times \text{Fix } \sigma & \xrightarrow{\text{enum}} & \text{List Patch} \\ \mathcal{O} \downarrow & & \downarrow \text{best} \\ \text{Fix}_a \sigma \times \text{Fix}_a \sigma & \xrightarrow{\text{tr}} & \text{Patch} \end{array}$$

Flag indicating copy:

```
data Fixa (s : Sum) : Set where
  ⟨·,·⟩ : Bool → [[s]]s (Fixa s) → Fixa s
```

Domain specific. The better the Oracle, the better the resulting patch.



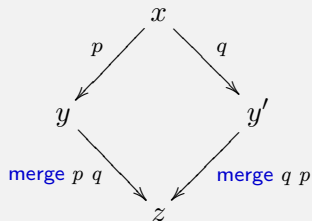
# Questions

- ▶ How can we represent data types?
- ▶ How can we represent patches on these data types?
- ▶ How can we compute such patches?
- ▶ **How can we merge such patches?**



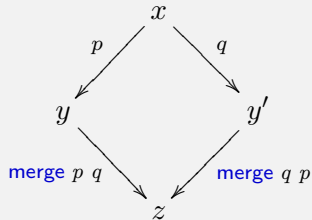
# Merging Square

Merging disjoint patches trivially commutes:



# Merging Square

Merging disjoint patches trivially commutes:



That is,

$$\text{apply} (\text{merge } p \ q) \circ \text{apply } p \equiv \text{apply} (\text{merge } q \ p) \circ \text{apply } q$$



## Related work

- ▶ There is a great deal of work on comparing (untyped) tree comparisons – but much less work that attempts to exploit the **type structure** that we have available.
- ▶ Lempink et al. & Vassena are a notable exception – but run a linear diff on the traversal of the tree. This it hard to guarantee that later operations – such as merging patches – produce well-formed trees.



# Looking ahead

## Work in Progress

- ▶ Prove properties about the efficient route of the “computing patches” square.
- ▶ Empirical Validation: Analysis of Clojure data from GitHub repos. Conducted by our MSc Giovanni Garufi.

## Future Work

- ▶ Implement a *proof-of-concept* in Haskell.
- ▶ Incorporate conflicts to our model.





# Type-directed Diffing of Structured Data

Victor Cacciari Miraldo, Pierre-Évariste Dagand  
and Wouter Swierstra

January 5, 2018

