

# An Efficient Algorithm for Type-Safe Structural Diffing

VICTOR CACCIARI MIRALDO, Utrecht University, The Netherlands

WOUTER SWIERSTRA, Utrecht University, The Netherlands

Effectively computing the difference between two version of a source file has become an indispensable part of software development. The *de facto* standard tool used by most version control systems is the UNIX `diff` utility, that compares two files on a line-by-line basis without any regard for the *structure* of the data stored in these files. This paper presents an alternative *datatype generic* algorithm for computing the difference between two values of *any* algebraic datatype. This algorithm maximizes sharing between the source and target trees, while still running in linear time. Finally, this paper demonstrates that by instantiating this algorithm to the Lua abstract syntax tree and mining the commit history of repositories found on GitHub, the resulting patches can often be merged automatically, even when existing technology has failed.

CCS Concepts: • **Software and its engineering** → **General programming languages**.

Additional Key Words and Phrases: Generic Programming, diff, Version Control, Haskell

## ACM Reference Format:

Victor Cacciari Miraldo and Wouter Swierstra. 2019. An Efficient Algorithm for Type-Safe Structural Diffing. *Proc. ACM Program. Lang.* 3, ICFP, Article 113 (August 2019), 29 pages. <https://doi.org/10.1145/3341717>

## 1 INTRODUCTION

The UNIX `diff` [Hunt and McIlroy 1976] is an essential tool in modern software development. It has seen a number of use cases ever since it was created and lies at the heart of today's Software Version Control Systems. Tools such as git, mercurial and darcs, that enable multiple developers to collaborate effectively, are all built around the UNIX `diff` utility, that computes a patch between two versions of a file. It compares files on a line-by-line basis attempting to share as many lines as possible between the source and the destination files.

A consequence of the *by line* granularity of the UNIX `diff` is its inability to identify more fine grained changes in the objects it compares. For example, if two parts of a program were changed, but happen to be printed on the same line, the UNIX `diff` sees this as a *single* change. Ideally, however, the objects under comparison should dictate the granularity of change to be considered. This is precisely the goal of *structural differencing* tools.

In this paper we present an efficient datatype-generic algorithm to compute the difference between two elements of any mutually recursive family. In particular, our algorithm readily works over the abstract syntax tree of a programming language—thereby enabling, for example, two changes that work on separate parts of the AST to be trivially merged, even if they appear to be on the same line in the source file. We have implemented our algorithm in Haskell and make heavy use of its datatype generic programming capabilities.

---

Authors' addresses: Victor Cacciari Miraldo, Information and Computing Sciences, Utrecht University, Princetonplein, 5, Utrecht, Utrecht, 3584 CC, The Netherlands, [v.cacciarimiraldo@uu.nl](mailto:v.cacciarimiraldo@uu.nl); Wouter Swierstra, Information and Computing Sciences, Utrecht University, Princetonplein, 5, Utrecht, Utrecht, 3584 CC, The Netherlands, [w.swierstra@uu.nl](mailto:w.swierstra@uu.nl).

---



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2019 Copyright held by the owner/author(s).

2475-1421/2019/8-ART113

<https://doi.org/10.1145/3341717>

In general, we aim to compute the difference between two values of type  $a$ , and represent these changes in some type,  $Patch\ a$ . The  $diff$  function computes these differences between two values of type  $a$ , and  $apply$  attempts to transform one value according to the information stored in the  $Patch$  provided to it.

$$\begin{aligned} diff &:: a \rightarrow a \rightarrow Patch\ a \\ apply &:: Patch\ a \rightarrow a \rightarrow Maybe\ a \end{aligned}$$

Note that the  $apply$  function may fail, for example, when attempting to delete data that is not present. Yet when it succeeds, the  $apply$  function must return a value of type  $a$ . This may seem like an obvious design choice, but this property does not hold for the approaches [Asenov et al. 2017; Falleri et al. 2014] using xml or json to represent their abstract syntax trees, where the result of applying a patch may produce ill-typed results, i.e., schema violations.

Naturally, not every definition of  $Patch$ ,  $diff$  and  $apply$  will solve our problem. We expect certain properties of our  $diff$  and  $apply$  functions. The first being *correctness*: the patch that  $diff\ x\ y$  computes can be used to faithfully reproduce  $y$  from  $x$ .

$$\forall x\ y. apply\ (diff\ x\ y)\ x \equiv Just\ y$$

The  $apply$  function is inherently partial and *correctness* only requires  $apply$  to succeed in one particular instance—but what should happen when applying a patch to a different value than the one used to create the input patch? We argue that the  $apply$  function should only fail when strictly necessary. In particular, if there are no changes, the patch should represent a *no-op*, and its application should be the identity:

$$\forall x\ y. apply\ (diff\ x\ x)\ y \equiv Just\ y$$

This captures the idea that a patch that does not make any modifications must be applicable to any value.

Finally, the last important properties stem from a practical perspective. We need both the  $diff$  and  $apply$  functions to be computationally efficient.

The UNIX  $diff$  [Hunt and McIlroy 1976] satisfies these properties for the specific type of lines of text, or,  $a \equiv [String]$ . It represents patches as a series of insertions, deletions and copies of lines and works by enumerating all possible patches that transform the source into the destination and chooses the ‘best’ such patch. There have been several attempts at generalizing these results to handle arbitrary datatypes [Lempsink et al. 2009; Miraldo et al. 2017], but following the same recipe: enumerate all combinations of insertions, deletions and copies that transform the source into the destination and choose the ‘best’ one. We argue that this design has two weaknesses when generalized to work over arbitrary types: (A) the non-deterministic nature of the design makes the algorithms inefficient, and (B), there exists no canonical ‘best’ patch and the choice is arbitrary.

We illustrate this last point with the example in Figure 1. The existing datatype generic approaches with insertions, deletions and copies typically perform a preorder traversal of the trees, copying over constructors whenever possible. Yet if we want to transform a binary tree  $Bin\ t\ u$  into  $Bin\ u\ t$  using only these operations, we will be forced to choose between copying  $t$  or  $u$ , but never both. The choice of which subtree to copy becomes arbitrary and unpredictable. To make matters worse, the non-determinism such choice points introduce makes algorithms intractably slow.

The central design decision underlying the UNIX  $diff$  tool is to *copy data whenever possible*. Yet this example shows that using only insertions, deletions and copies *limits* the opportunities for copying data. In the presence of richly structured data beyond lines of text, this becomes especially problematic.

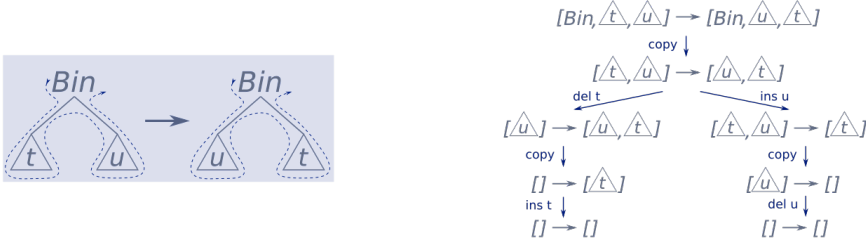


Fig. 1. Visualization of a *diff* (*Bin t u*) (*Bin u t*) using insertions, deletions and copies only

This paper explores a novel direction for differencing algorithms: rather than restricting ourselves to insertions, deletions, and copy operations, we allow the arbitrary reordering, duplication, and contraction of subtrees. Not only does this restrict the inherent non-determinism in the problem, making it *easier* to compute patches, this also *increases* the opportunities for copying. More specifically, this paper makes the following novel contributions:

- This paper defines a datatype generic *diff* function that computes a patch between two algebraic datatypes that is efficient in both time and space. This *diff* function supports the duplication and permutation of subtrees, and satisfies all the desired properties outlined above. We illustrate this algorithm by first defining a specific instance (Section 2), then presenting a generic version capable of handling arbitrary mutually recursive families of datatypes (Section 3).
- Initially, we present our diff algorithm assuming the existence of an oracle capable of detecting all possible copying opportunities. We give a practical, generic, implementation of this oracle that is correct modulo cryptographic hash collisions and runs in amortized constant time (Section 3.3).
- We show how the representation for patches used in this paper enables *disjoint* patches to be *merged* automatically (Section 4).
- Finally, we have instantiated our algorithm to the abstract syntax tree of Lua and collected historical data regarding merge conflicts from popular GitHub repositories. We show how our naive merging algorithm is already capable of resolving more than 10% of the merge conflicts encountered automatically, while still offering competitive performance (Section 5).

## 2 TREE DIFFING: A CONCRETE EXAMPLE

Before exploring the generic implementation of our algorithm, let us look at a simple, concrete instance first. This example sets the stage for the the generic implementation that follows (Section 3.2). Throughout this section we will explore the central ideas from our algorithm instantiated for the type of 2-3-trees:

```
data Tree23 = Leaf
            | Node2 Tree23 Tree23
            | Node3 Tree23 Tree23 Tree23
```

The central concept of our work is the encoding of a *change*. Unlike previous work [Klein 1998; Lempsink et al. 2009; Miraldo et al. 2017] which is based on tree-edit-distance [Bille 2005] and hence, uses only insertions, deletions and copies of the constructors encountered during the preorder traversal of a tree (Figure 1), we go a step further. We explicitly model permutations, duplications and contractions of subtrees within our notion of *change*. Where contraction here denotes the

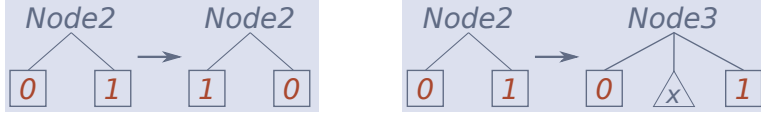


Fig. 2. Visualization of  $\text{diff}(\text{Node2 } t \ u) (\text{Node2 } u \ t)$  to the left and  $\text{diff}(\text{Node2 } t \ u) (\text{Node3 } t \ x \ u)$  on the right. Metavariables are shown inside a square. Entire trees are shown inside a triangle.

partial inverse of a duplication. The representation of a *change* between two values of type *Tree23*, then, is given by identifying the bits and pieces that must be copied from source to destination making use of permutations and duplications where necessary.

A new datatype, *Tree23C*  $\varphi$ , enables us to annotate a value of *Tree23* with holes of type  $\varphi$ . Therefore, *Tree23C MetaVar* represents the type of *Tree23* with holes carrying metavariables. These metavariables correspond to arbitrary trees that are *common subtrees* of both the source and destination of the change. These are exactly the bits that are being copied from the source to the destination tree. We refer to a value of *Tree23C* as a *context*. For now, the metavariables will be simple *Int* values but later on we will need to carry additional information.

**type** *MetaVar* = *Int*

**data** *Tree23C*  $\varphi$  = *Hole*  $\varphi$   
     | *LeafC*  
     | *Node2C* *Tree23C* *Tree23C*  
     | *Node3C* *Tree23C* *Tree23C* *Tree23C*

A *change* in this setting is a pair of such contexts. The first context defines a pattern that binds some metavariables, called the deletion context; the second, called the insertion context, corresponds to the tree annotated with the metavariables that are supposed to be instantiated by the bindings given by the deletion context.

**type** *Change23*  $\varphi$  = (*Tree23C*  $\varphi$ , *Tree23C*  $\varphi$ )

The change that transforms *Node2*  $t \ u$  into *Node2*  $u \ t$  is then represented by a pair of *Tree23C*, (*Node2C* (*Hole* 0) (*Hole* 1), *Node2C* (*Hole* 1) (*Hole* 0)), as seen in Figure 2. This change works on any tree built using the *Node2* constructor and swaps the children of the root. Note that it is impossible to define such swap operations in terms of insertions and deletions—as used by most diff algorithms.

## 2.1 Applying Changes

Applying a change is done by instantiating the metavariables in the deletion context and the insertion context:

*applyChange* :: *Change23* *MetaVar*  $\rightarrow$  *Tree23*  $\rightarrow$  *Maybe* *Tree23*  
*applyChange* ( $d, i$ )  $x$  = *del*  $d \ x$   $\gg$  *ins*  $i$

Naturally, if the term  $x$  and the deletion context  $d$  are *incompatible*, this operation will fail. Contrary to regular pattern-matching we allow variables to appear more than once on both the deletion and insertion contexts. Their semantics are dual: duplicate variables in the deletion context must match equal trees, and are referred to as *contractions*, whereas duplicate variables in the insertion context will duplicate trees. We use an auxiliary function within the definition of *del* to

make this check easier to perform. Given a deletion context  $ctx$  and source  $tree$ , the  $del$  function tries to associate all the metavariables in the context with a subtree of the input  $tree$ .

```
del :: Tree23C MetaVar → Tree23 → Maybe (Map MetaVar Tree23)
del ctx tree = go ctx tree empty
```

The  $go$  function, defined below, closely follows the structure of trees and contexts. Only when we reach a *Hole*, do we check whether we have already instantiated the metavariable stored there or not. If we have encountered this metavariable before, we check that both occurrences of the metavariable correspond to the same tree; if this is the first time we have encountered this metavariable, we simply instantiate the metavariable with the current tree.

```
go :: Tree23C → Tree23 → Map MetaVar Tree23 → Maybe (Map MetaVar Tree23)
go LeafC      Leaf      m = return m
go (Node2C x y) (Node2 a b) m = go x a m >> go y b
go (Node3C x y z) (Node3 a b c) m = go x a m >> go y b >> go z c
go (Hole i)      t        m = case lookup i m of
                                Nothing → return (M.insert i t m)
                                Just t'  → guard (t ≡ t') >> return m
go _             _        m = Nothing
```

We will refer to the result of  $del\ ctx\ tree$  as the *valuation* that instantiates the metavariables of  $ctx$  with subtrees of  $tree$ . Once we have obtained a such valuation, we substitute the variables in the insertion context with their respective values, to obtain the final tree. This phase fails when the change contains unbound variables. The  $ins$  function is defined below.

```
ins :: Tree23C MetaVar → Map MetaVar Tree23 → Maybe Tree23
ins LeafC      m = return Leaf
ins (Node2C x y) m = Node2 <$> ins x m <*> ins y m
ins (Node3C x y z) m = Node3 <$> ins x m <*> ins y m <*> ins z m
ins (Hole i)      m = lookup i m
```

## 2.2 Computing Changes

Next, we explore how to produce a change from a source and a destination, defining a  $changeTree23$  function. Intuitively, this function will try to exploit as many copy opportunities as possible. For now, we delegate the decision of whether a subtree should be copied or not to an oracle: assume we have access a function  $wcs :: Tree23 → Tree23 → Tree23 → Maybe MetaVar$ , short for “*which common subtree*”. The call  $wcs\ s\ d\ x$  returns *Nothing* when  $x$  is *not* a subtree of  $s$  and  $d$ ; if  $x$  is a subtree of both  $s$  and  $d$ , it returns *Just i*, for some metavariable  $i$ . The only condition we impose is injectivity of  $wcs\ s\ d$ : that is, if  $wcs\ s\ d\ x ≡ wcs\ s\ d\ y ≡ Just\ j$ , then  $x ≡ y$ . In other words, equal metavariables correspond to equal subtrees.

There is an obvious inefficient implementation for  $wcs$ , that traverses both trees searching for shared subtrees—hence postulating the existence of such an oracle is not a particularly strong assumption to make. In Section 2.4, we provide an efficient implementation for  $Tree23$ . For now, assuming the oracle exists allows for a clear separation of concerns. The  $changeTree23$  function merely has to compute the deletion and insertion contexts, using said oracle—the inner workings of the oracle are abstracted away cleanly.

```
changeTree23 :: Tree23 → Tree23 → Change23 MetaVar
changeTree23 s d = (extract (wcs s d) s, extract (wcs s d) d)
```

$$\begin{aligned}
a &= \text{Node2} (\text{Node2 } t \ k) \ u & \text{extract } (\text{wcs } a \ b) \ a &= \text{Node2C} (\text{Hole } 0) \ u \\
b &= \text{Node2} (\text{Node2 } t \ k) \ t & \text{extract } (\text{wcs } a \ b) \ b &= \text{Node2C} (\text{Hole } 0) (\text{Hole } 1) \\
\\
\text{postprocess } a \ b \ (\text{extract } (\text{wcs } a \ b) \ a) \ (\text{extract } (\text{wcs } a \ b) \ b) \\
&= (\text{Node2C} (\text{Hole } 0) \ u, \text{Node2C} (\text{Hole } 0) \ t)
\end{aligned}$$

Fig. 3. Example of erroneous context extraction due to nested common subtrees

The *extract* function receives an oracle and a tree. It traverses its argument tree, looking for opportunities to copy subtrees. It repeatedly consults the oracle, to determine whether or not the current subtree should be shared across the source and destination. If that is the case, we want our change to *copy* such subtree. That is, we return a *Hole* whenever the second argument of *extract* is a common subtree according to the oracle. If the oracle returns *Nothing*, we move the topmost constructor to the context being computed and recurse over the remaining subtrees.

```

extract :: (Tree23 → Maybe MetaVar) → Tree23 → Tree23C MetaVar
extract o t = maybe (peel t) Hole (o t)
  where peel Leaf      = LeafC
        peel (Node2 a b) = Node2C (extract o a) (extract o b)
        peel (Node3 a b c) = Node3C (extract o a) (extract o b) (extract o c)

```

Note that had we used a version of *wcs* that only returns a boolean value we would not know what metavariable to use when a subtree is shared. Returning a value that uniquely identifies a subtree allows us to keep the *extract* function linear in the number of constructors in *x* (disregarding the calls to our oracle for the moment).

This iteration of the *changeTree23* function has a subtle bug: not all common subtrees can be copied. In particular, we cannot copy a tree *t* that occurs as a subtree of the source and destination, but also appears as a subtree of another, larger common subtree. One such example is shown in Figure 3, where the oracle claims that both *Node2 t k* and *t* are common subtrees. As *t* also occurs by itself one of the extracted contexts will contain an unbound metavariable. This will trigger an error when trying to apply the corresponding change. In this example, applying the change from Figure 3 would trigger such error when the *ins* function branch for the *Hole* constructor and attempts to lookup the tree associated with metavariable 1.

One way to solve this is to introduce an additional postprocessing step that substitutes the variables that occur exclusively in the deletion or insertion context by their corresponding tree. We can implement this postprocessing step using two calls to the *del* function we saw previously: one for the deletion context against the source tree and another for the insertion context against the destination tree. The resulting information is then used to replace *unused* or *undeclared* metavariables with the tree to which they correspond. We postpone the implementation until its generic incarnation in Section 3.2.

```

postprocess :: Tree23 → Tree23 → Tree23C MetaVar → Tree23C MetaVar
              → (Tree23C MetaVar, Tree23C MetaVar)

```

We fix the previous *changeTree23* by postprocessing the extracted contexts. The new version of *changeTree23* will only produce closed changes, where each deletion and insertion context have the same set of metavariables. Intuitively, this means that every variable that is declared is used and vice-versa.

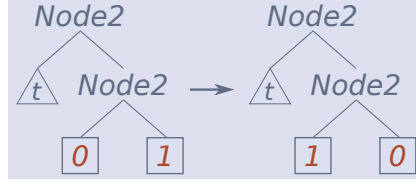


Fig. 4. A *Change23* with a redundant *Node2* *t* in both the deletion and insertion contexts.

$changeTree23 :: Tree23 \rightarrow Tree23 \rightarrow Change23\ MetaVar$

$changeTree23\ s\ d = postprocess\ s\ d\ (extract\ (wcs\ s\ d)\ s)\ (extract\ (wcs\ s\ d)\ d)$

Assuming that  $wcs\ s\ d$  correctly assigns metavariables to *all* common subtrees of  $s$  and  $d$ , it is not hard to see that our implementation already satisfies the specification we formulated in the introduction:

**Correctness** Assuming  $wcs$  is correct,

$$\forall x\ y. applyTree23\ (changeTree23\ x\ y)\ x \equiv Just\ y$$

**Preciseness** Assuming  $wcs$  is correct,

$$\forall x\ y. applyTree23\ (changeTree23\ x\ x)\ y \equiv Just\ y$$

**Time Efficiency** On the worst case, we perform one query to the oracle per constructor in our trees. With a amortized constant time  $wcs$ , Section 3.3, our algorithm is linear on the number of constructors in the source and destination trees.

**Space Efficiency** The size of a *Change23 MetaVar* is, on average, smaller than storing its source and destination tree completely. On the worst case, where there is no common subtree, they have the same size. This is also true of the Unix `diff` utility, when comparing two files that do not share a single line of text.

Although correct with respect to our specification, there is still room for improvement. A call to  $changeTree23\ x\ y$  yields a *single Change23*, consisting of a pair of insertion and deletion contexts. When  $x$  and  $y$  resemble one another these contexts may store a great deal of redundant information as many constructors appearing in both contexts will be ‘deleted’, and then ‘inserted’, as shown in Figure 4. More importantly, however, is the fact that when looking at the deletion or insertion contexts of Figure 4, we do not know whether a constructor is being copied over or not, which hinders our capacity to easily merge patches (Section 4). For example, while merging a patch if we see a change that looks like *Change (Hole i) (Hole i)*, we would like to merge it immediately.

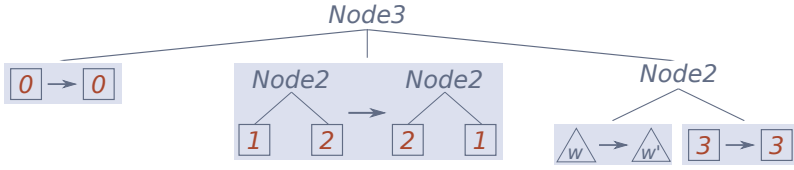
We must be careful to make sure that the domain of the patch after minimizing changes is superset of the original.

### 2.3 Minimizing Changes: Computing Patches

The process of minimizing and isolating the changes starts by identifying the redundant part of the contexts. That is, the constructors that show up as a prefix in *both* the deletion and the insertion context. They are essentially being copied over and we want to make this fact explicit by separating them into what we call the *spine* of the patch. This step will help us reason about patches later on. If a constructor is in the *spine*, we know it has been copied, if it shows up in a change, we know it was either deleted or inserted. The spine will then contain changes—pairs of an insertion and deletion context—in its leaves:

$type\ Patch23 = Tree23C\ (Change23\ MetaVar)$





$Node3C$  ( $Hole$  ( $Hole$  0,  $Hole$  0))  
 ( $Hole$  ( $Node2C$  ( $Hole$  0) ( $Hole$  1),  $Node2C$  ( $Hole$  1) ( $Hole$  0))  
 ( $Node2C$  ( $Hole$  ( $tree23toC$   $w$ ,  $tree23toC$   $w'$ ))  
 ( $Hole$  ( $Hole$  3,  $Hole$  3)))

Fig. 5. Graphical and textual representation of the patch that transforms the value  $Node3\ t\ (Node2\ u\ v)\ (Node2\ w\ x)$  into the value  $Node3\ t\ (Node2\ v\ u)\ (Node2\ w'\ x)$ . The  $tree23toC$  function converts a  $Tree23$  into a  $Tree23C$  in the canonical way.

A *patch* consists in a spine with *changes* inside of it. Figure 5 illustrates a value of type  $Patch23$ , where the *changes* are visualized with a shaded background in the leaves of the spine. Note that the changes contains only the necessary constructors to make sure that all metavariables that are used in the insertion context are defined in the deletion context. This keeps changes small and isolated, making them easier to merge.

In this section we will discuss how to take the results of  $changeTree23$  and transform them into a  $Patch23$ . The first step to compute a patch from a change is identifying its *spine*. That is, the constructors that are present in both the deletion and insertion contexts. We are essentially splitting a monolithic change into the *greatest common prefix* of the insertion and deletion contexts, leaving smaller changes on the leaves of this prefix:

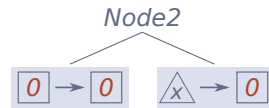
$gcp :: Tree23C\ var \rightarrow Tree23C\ var \rightarrow Tree23C\ (Change23\ var)$   
 $gcp\ LeafC\ LeafC = LeafC$   
 $gcp\ (Node2C\ a1\ b1)\ (Node2C\ a2\ b2) = Node2C\ (gcp\ a1\ a2)\ (gcp\ b1\ b2)$   
 $gcp\ (Node3C\ a1\ b1\ c1)\ (Node3C\ a1\ b2\ c2) = Node3C\ (gcp\ a1\ a2)\ (gcp\ b1\ b2)\ (gcp\ c1\ c2)$   
 $gcp\ a\ b = Hole\ (a, b)$

In the last case of the  $gcp$  function either  $a$  and  $b$  are both holes or the constructors disagree, hence they do *not* belong in the common prefix.

One might be tempted to take the results of  $changeTree23C$ , pipe them into the  $gcp$  function directly. Yet, the *greatest common prefix* consumes all the possible constructors leading to disagreeing parts of the contexts where this might be too greedy. We must be careful not to break bindings as shown below:

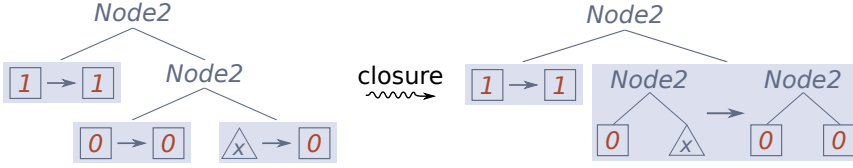
$--\ prob = changeTree23\ (Node2\ t\ t)\ (Node2\ x\ t)$   
 $prob :: Change23\ MetaVar$   
 $prob = Change\ (Node2C\ (Hole\ 0)\ (Hole\ 0))$   
 $\quad\quad\quad ,\ Node2C\ x\ (Hole\ 0))$

$gcp\ prob \equiv$



In this example, the second change contains a  $Hole\ 0$  that does not occur in the deletion context, and is hence *unbound*. To address this problem, we go over the result from our call to  $gcp$ , pulling changes up the tree until each change is closed, that is, the set of variables in both contexts is identical. We call this process the *closure* of a patch and declare a function to compute this below.



Fig. 6. Graphical representation of the *closure* function.

We have illustrated the process of *closure* in Figure 6. Note that in both the input and output for the *closure* function the subtree  $x$  appears on the deletion context. Moreover, the *closure* functions only bubbles up the minimal number of constructors to ensure all changes are closed. In particular, the *Node2* constructor at the root is still part of the spine after the call to *closure*.

$\text{closure} :: \text{Tree23C } (\text{Change23 MetaVar}) \rightarrow \text{Patch23}$

Although the *closure* function apparently always returns a patch, its implementation might fail if there exists no way of closing all the changes. In our case, this will never happen as we know that *changeTree23* outputs a closed change. In the worst case, the resulting spine will be empty—but the change will certainly be closed. That is due to *postprocess*, which could have been merged with *closure*. We kept them separate for clarity. Section 3.2 comes back to the *closure* function in more detail on its generic incarnation. For now, it is more important to understand that it facilitates merging our patches, later on (Section 4). It is worth mentioning that the result of *closure p* should be a patch that can be applied to *at least as many elements as p*. In some corner cases *closure* might enlarge the domain of a patch by breaking some contractions. This is not an issue, however, as being able to apply a patch to more elements is a good thing.

As soon as every change within the spine has been closed, we have a *patch*. The final *diff* function for *Tree23* is then defined as follows:

$\text{diffTree23} :: \text{Tree23} \rightarrow \text{Tree23} \rightarrow \text{Patch23}$

$\text{diffTree23 } s \ d = \text{closure } \$ \text{ gcp } \$ \text{ changeTree23 } s \ d$

We could now define the *applyPatch23* function that applies a *patch*, rather than the *applyChange23* we saw previously. This is done by traversing the object tree and the spine of the patch until a change is found and applying that change to the *local* subtrees in question.

## 2.4 Defining the Oracle for *Tree23*

In order to have a working version of our diff algorithm for *Tree23* we must provide the *wcs* implementation. Recall that the *wcs* function, *which common subtree*, has type  $\text{Tree23} \rightarrow \text{Tree23} \rightarrow \text{Tree23} \rightarrow \text{Maybe MetaVar}$ . Given a fixed  $s$  and  $d$ ,  $\text{wcs } s \ d \ x$  returns *Just i* if  $x$  is the  $i^{\text{th}}$  subtree of  $s$  and  $d$  and *Nothing* if  $x$  does not appear in  $s$  or  $d$ . One implementation of this function computes the intersection of all the subtrees in  $s$  and  $d$ , and then search for the subtree  $x$  the resulting list. Enumerating all the subtrees of any *Tree23* is straightforward:

$\text{subtrees} :: \text{Tree23} \rightarrow [\text{Tree23}]$

$\text{subtrees Leaf} = [\text{Leaf}]$

$\text{subtrees (Node2 } x \ y) = \text{Node2 } x \ y : (\text{subtrees } x \# \text{ subtrees } y)$

$\text{subtrees (Node3 } x \ y \ z) = \text{Node3 } x \ y \ z : (\text{subtrees } x \# \text{ subtrees } y \# \text{ subtrees } z)$

It is now straightforward to implement the *wcs* function: we compute the intersection of all the subtrees of  $s$  and  $d$  and use this list to determine whether the argument tree occurs in both  $s$  and  $d$ .

This check is done with *elemIndex* which returns the index of the element, when it occurs in the list.

```
wcs :: Tree23 → Tree23 → Tree23 → Maybe MetaVar
wcs s d x = elemIndex x (subtrees s ∩ sutrees d)
```

This implementation, however, is not particularly efficient. The inefficiency comes from two places: firstly, checking trees for equality is linear in the size of the tree; furthermore, enumerating all subtrees is exponential. If we want our algorithm to be efficient we *must* have an amortized constant-time *wcs*.

To tackle the first issue and efficiently compare trees for equality we will be using cryptographic hash functions [Menezes A. J. and Vanstone [n. d.]] to construct a fixed length bitstring that uniquely identifies a tree modulo hash collisions. Said identifier will be the hash of the root of the tree, which will depend on the hash of every subtree, much like a *merkle tree* [Merkle 1988]. Suppose we have a function *merkleRoot* that computes some suitable identifier for every tree, we can compare trees efficiently by comparing their associated identifiers:

**instance** *Eq Tree23 where*

```
t ≡ u = merkleRoot t ≡ merkleRoot u
```

The definition of *merkleRoot* function is straightforward. It is important that we use the *merkleRoot* of the parts of a *Tree23* to compute the *merkleRoot* of the whole. This construction, when coupled with a cryptographic hash function, call it *hash*, is what guarantee injectivity modulo hash collisions.

```
merkleRoot :: Tree23 → Digest
merkleRoot Leaf = emptyDigest
merkleRoot (Node2 x y) = hash (concat ["node2", merkleRoot x, merkleRoot y])
merkleRoot (Node3 x y z) = hash (concat ["node3", merkleRoot x, merkleRoot y, merkleRoot z])
```

Note that although it is theoretically possible to have false positives, when using a cryptographic hash function the chance of collision is negligible and hence, in practice, they never happen [Menezes A. J. and Vanstone [n. d.]]. Nonetheless, it is easy to detect when a collision has occurred in our algorithm; consequently, we chose to ignore this issue.

Recall we are striving for a constant time ( $\equiv$ ) implementation, but the ( $\equiv$ ) definition above is still linear, we recompute the hash on every comparison. We fix this by caching the hash associated with every node of a *Tree23*. This is done by the *decorate* function, Figure 7.

**data** *Tree23H = LeafH*

```
| Node2H (Tree23H, Digest) (Tree23H, Digest)
| Node3H (Tree23H, Digest) (Tree23H, Digest) (Tree23H, Digest)
```

```
decorate :: Tree23 → Tree23H
```

We omit the implementation of *decorate* for brevity even if it is straightforward. Moreover, a generic version is introduced in Section 3.3. This enables us to define a constant time *merkleRoot* function, shown below, which makes the ( $\equiv$ ) function run in constant time.

```
merkleRoot :: Tree23H → Digest
merkleRoot LeafH = emptyDigest
merkleRoot (Node2H (_, hx) (_, hy)) = hash (encode "2" + hx + hy)
merkleRoot (Node3H (_, hx) (_, hy) (_, hz)) = hash (encode "3" + hx + hy + hz)
```

The second source of inefficiency, enumerating all possible subtrees, can be addressed by choosing a better data structure. In order to check whether a tree *x* is a subtree of a fixed *s* and *d*, it suffices

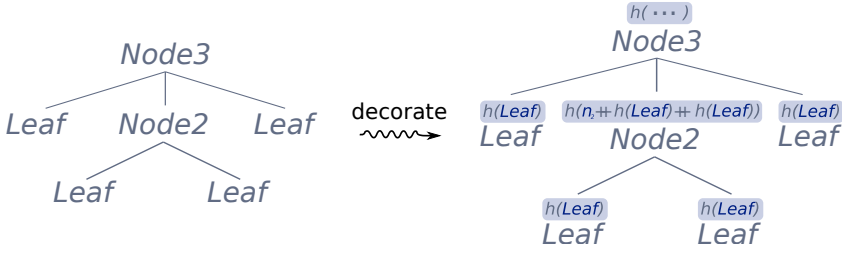


Fig. 7. Example of a merkelized *Tree23*, where  $n_2$  is some fixed identifier and  $h$  is a hash function.

to check whether the merkle root of  $x$  appears in a “database” of the common merkle roots of  $s$  and  $d$ . Given that a *Digest* is just a *[Word]*, the optimal choice for such “database” is a Trie [Brass 2008] mapping a *[Word]* to a *MetaVar*. Trie lookups are efficient and hardly depend on the number of elements in the trie. In fact, our lookups run in amortized constant time here, as the length of a *Digest* is fixed.

Finally, we are able to write our efficient *wcs* oracle that concludes the implementation of our algorithm for the concrete *Tree23* type. The *wcs* oracle will now receive *Tree23H*, i.e., trees annotated with their merkle roots at every node, and will populate the “database” of common digests.

$wcs :: Tree23H \rightarrow Tree23H \rightarrow Tree23H \rightarrow Maybe\ MetaVar$

$wcs\ s\ d = lookup\ (mkTrie\ s \cap mkTrie\ d) \circ merkleRoot$

where

$(\cap) :: Trie\ k\ v \rightarrow Trie\ k\ u \rightarrow Trie\ k\ v$

$lookup :: Trie\ k\ v \rightarrow [k] \rightarrow Maybe\ v$

$mkTrie :: Tree23H \rightarrow Trie\ Word\ MetaVar$

The use of cryptographic hashes is unsurprising. They are almost folklore for speeding up a variety of computations. It is important to notice that the efficiency of our algorithm comes from our novel representation of patches combined with a amortized constant time *wcs* function. Without being able to duplicate or permute subtrees, the algorithm would have to backtrack in a number of situations.

### 3 TREE DIFFING GENERICALLY

In Section 2 we provided a simple algorithm to solve the differencing problem for 2-3-trees. We began by creating the type of contexts over *Tree23*, which consisted of annotating a *Tree23* with a *metavariable* constructor. Later, assuming the existence of an oracle that determines whether or not an arbitrary tree is a subtree of the source and the destination, we compute a value of type *Change23 MetaVar* from a *Tree23*. How to compute a *Patch23* given a *Change23* by *minimizing* the changes and isolating them in the *spine*. In this section we show how can we write that same algorithm in a generic fashion, working over any mutually recursive family of datatypes. The code in this section generalizes the example from the previous section, but we assume some familiarity with generic programming in modern Haskell. Readers unfamiliar with these concepts may safely skip this section on a first reading.

### 3.1 Background on Generic Programming

Firstly, let us briefly review the `generics-mrsop` [Miraldo and Serrano 2018] library, that we will use to define a generic version of our algorithm. This library follows the *sums-of-products* school of generic programming [de Vries and Löh 2014] and, additionally, enables us to work with mutually recursive families. This is particularly important for this algorithm, as the abstract syntax trees of many programming languages consist of mutually recursive types for expressions, statements, methods, classes and other language constructs.

Take the `Tree23` type from Section 2. Its structure can be seen in a *sum-of-products* fashion through the `Tree23SOP` type given below. It represents the shape in which every `Tree23` comes and consists in two nested lists of *atoms*. The outer list represents the choice of constructor, and packages the *sum* part of the datatype whereas the inner list represents the *product* of the fields of a given constructor. The `'·` notation represents a value that has been promoted to the type level [Yorgey et al. 2012].

```
type Tree23SOP = '['[]
                , '[I 0, I 0]
                , '[I 0, I 0, I 0]]
```

The atoms, in this case only `I 0`, represent a recursive position referencing the first type in the family. In fact, a mutually recursive family is described by a *list of sums-of-products*: one for each element in the family. We overload the word “code” in singular or plural to mean the representation of a datatype, or the representation of a family, respectively. The context should make clear the distinction. For example, `Tree23SOP` is the code of type `Tree23` and `Tree23Code` is the codes for the mutually recursive family, which in this case, contains only one type.

```
type Tree23Code = '['Tree23SOP]
```

The `generics-mrsop` uses the type `Atom` to distinguish whether a field is a recursive position referencing the  $n$ -th type in the family, `I n`, or a opaque type, for example, `Int` or `Bool`, which are represented by `K KInt`, `K KBool`.

Let us now take a mutually recursive family with more than one element and see how it is represented within the `generics-mrsop` framework. The mutually recursive family containing types `Zig` and `Zag` has its codes defined as a list of codes, one for each member of the family:

```
data Zig = Zig Int | ZigZag Zag      type ZigCodes = '['['[K KInt] , '[I 1]]
data Zag = Zag Bool | ZagZig Zig    , '['[K KBool], '[I 0]]
                                   ]
```

Note that the codes come in the same order as the elements of the family. The code for `Zig` is the first element of the `ZigCodes` type level list. It consists in two lists, since `Zig` has two constructors. One receives a value of type `Int`, the other consists in a recursive call to the second element of the family. The code acts as a recipe that the *representation* functor must follow in order to interpret those into a type of kind `*`.

The representation is defined by the means of  $n$ -ary sums (*NS*) and products (*NP*) that work by induction on the *codes* and one interpretation for atoms (*NA*). Their definition together with their respective elimination principles can be found in Figure 8.

The *NS* type is responsible for determining the choice of constructor whereas the *NP* applies a representation functor to all the fields of the selected constructor. Finally, *NA* distinguishes between representation of a recursive position from an opaque type. Although the `generics-mrsop` provides a way to customize the set of opaque types used, this is not central to the developments in this

|  |  |
|--|--|
| <pre> <b>data</b> <i>NS</i> :: (k → *) → [k] → * <b>where</b>   <i>Here</i> :: f x      → <i>NS</i> f (x ' : xs)   <i>There</i> :: <i>NS</i> f xs → <i>NS</i> f (x ' : xs)  <b>data</b> <i>NP</i> :: (k → *) → [k] → * <b>where</b>   <i>Nil</i>  :: <i>NP</i> f '[]   <i>Cons</i> :: f x → <i>NP</i> f xs → <i>NP</i> f (x ' : xs)  <b>data</b> <i>NA</i> :: (Nat → *) → Atom → * <b>where</b>   <i>NA<sub>I</sub></i> :: φ i      → <i>NA</i> φ (I i)   <i>NA<sub>K</sub></i> :: <i>Opq</i> k → <i>NA</i> φ (K k) </pre> | <pre> <i>elimNS</i> :: (∀ at . f at → a) → <i>NS</i> f s → a <i>elimNS</i> f (<i>There</i> s) = <i>elimNS</i> f s <i>elimNS</i> f (<i>Here</i> x) = f x  <i>elimNP</i> :: (∀ at . f at → a) → <i>NP</i> f p → [a] <i>elimNP</i> f <i>Nil</i>           = [] <i>elimNP</i> f (<i>Cons</i> x xs) = f x : <i>elimNP</i> f xs  <i>elimNA</i> :: (∀ ix . f x → a) → (∀ k . g k → a)           → <i>NA</i> f g at → a <i>elimNA</i> f g (<i>NA<sub>I</sub></i> x) = f x <i>elimNA</i> f g (<i>NA<sub>K</sub></i> x) = g x </pre> |
|--|--|

Fig. 8. Interpretation and elimination principles for each component of a sum-of-products code.

paper and hence we will assume a type *Opq* that interprets the necessary atom types, i.e., *Int*, *Bool*, etc. We refer the interested reader to the original paper [Miraldo and Serrano 2018] for more information. Nevertheless, we define the representation functor *Rep* as the composition of the interpretations of the different pieces:

**type** *Rep* φ = *NS* (*NP* (*NA* φ))

Finally, we tie the recursive knot with a functor of kind *Nat* → \* that is passed as a parameter to *NA* in order to interpret the recursive positions. The familiar reader might recognize it as the indexed least fixpoint:

**newtype** *Fix* (codes :: '[ '[Atom] ]]) (ix :: Nat)  
 = *Fix* { *unFix* :: *Rep* (*Fix* codes) (*Lkup* codes ix) }

Here, *Lkup* codes ix denotes the type level lookup of the element with index ix within the list codes. This type family throws a type error if the index is out of bounds. The generic versions of the constructors of type *Zig* are given by:

```

gzig :: Int → Fix ZigCodes 0
gzig n = Fix (Here (Cons (NAK (OpqInt n)) Nil))

gzigzag :: Fix ZigCodes 1 → Fix ZigCodes 0
gzigzag zag = Fix (There (Here (Cons (NAI zag) Nil)))

```

One of the main benefits of the *sums-of-products* approach to generic programming is that it enables us to pattern match generically. In fact, we can state that a value of a type consists precisely of a choice of constructor and a product of its fields by defining a *view* type. For example, we take the *constructor* of a generic type to be:

```

data Constr :: [[k]] → Nat → * where
  CZ :: Constr (x ' : xs) Z
  CS :: Constr xs c → Constr (x ' : xs) (S c)

```

The *Constr* sum c type represents a predicate indicating that c is a valid constructor for sum, that is, it is a valid index into the type level list sum. This enables us to define a *View* over a value of a sum type to be a choice of constructor and corresponding product. We can then unwrap a *Fix* codes i value into its topmost constructor and a product of its fields with the *sop* function.

```

data View :: (Nat → *) → '[Atom] → * where
  Tag :: Constr sum c → NP (NA φ) (Lkup sum c) → View φ sum
  sop :: Fix codes i → View (Fix codes) (Lkup codes i)

```

This brief introduction covers the basics of generic programming in Haskell that we will use in this paper. We refer the interested reader to the literature [de Vries and Löh 2014; Miraldo and Serrano 2018] for a more thorough overview.

### 3.2 Representing and Computing Changes, Generically

Section 2 presented one particular instance of our differencing algorithm. In this section, we will generalize the definition using the generic programming techniques we have just seen.

We start defining the generic notion of context, called *Tx*. Analogously to *Tree23C* (Section 2), *Tx* enables us to augment mutually recursive family with type holes. This type construction is crucial to the representation of patches. This construction can be done for any mutually recursive family.

We can read *Tx codes φ at* as the element of the mutually recursive family *codes* indexed by *at* augmented with holes of type *φ*. Its definition follows:

```

data Tx :: [[[Atom]]] → (Atom → *) → Atom → * where
  TxHole :: φ at → Tx codes φ at
  TxOpq :: Opq k → Tx codes φ (K k)
  TxPeel :: Constr (Lkup codes i) c
           → NP (Tx codes φ) (Lkup (Lkup codes i) c)
           → Tx codes φ (I i)

```

Looking at the definition of *Tx*, we see that its values consist in either a *typed* hole, some opaque value, or a constructor and a product of fields. The *TxPeel* follows very closely the *View* type from Section 3.1. The *Tx* type is, in fact, the indexed free monad over the *Rep*. The return method is just *TxHole*, and the multiplication is given by:

```

txJoin :: Tx codes (Tx codes φ) at → Tx codes φ at
txJoin (TxHole tx) = tx
txJoin (TxOpq opq) = TxOpq opq
txJoin (TxPeel c d) = TxPeel c (mapNP txJoin p)

```

Essentially, a value of type *Tx codes φ at* is a value of type *NA (Fix codes) at* augmented with holes of type *φ*. To illustrate this, let us define the injection of a *NA* into a *Tx*:

```

na2tx :: NA (Fix codes) at → Tx codes φ at
na2tx (NAK k) = TxOpq k
na2tx (NAI (Fix x)) = case sop x of Tag c p → TxPeel c (mapNP na2tx p)

```

The inverse operations is partial. We can translate a *Tx* into an *NA* when the *Tx* has no holes:

```

tx2na :: Tx codes φ at → Maybe (NA (Fix codes) at)
tx2na (TxHole _) = Nothing
tx2na (TxOpq k) = NAK k
tx2na (TxPeel c txs) = inj c <$> mapNPM tx2na txs

```

*Generic Representation of Changes.* With a generic notion of contexts, we can go ahead and define our generic *Change* type. We use a pair of *Tx* exactly as in Section 2: one deletion context and one

insertion context. This time, however, we define a new datatype to be able to partially apply its type arguments later on.

**data** *Change* codes  $\varphi$  at = *Change* (*Tx* codes  $\varphi$  at) (*Tx* codes  $\varphi$  at)

The interpretation for the metavariables, *MetaVar*, now carries the integer representing the metavariable itself but also carries information to identify whether this metavariable is supposed to be instantiated by a recursive member of the family or a opaque type. We do so by carrying a singleton [Eisenberg and Weirich 2012] of type *NA*. This enables the compiler to gain knowledge over *at* when pattern-matching, which is important purely from the generic programming perspective. Without this information we would not be able to write a well-typed application function, for instance. We must know the types of the values supposed to be matched against a metavariable to ensure we will produce well-typed trees.

**data** *MetaVar* at = *MetaVar* Int (*NA* (*Const Unit*) at)

The type of changes over *Tree23* can now be written using the generic representation for changes and metavariables.

**type** *ChangeTree23* = *Change Tree23Code MetaVar* (*I 0*)

We can read the type above as the type of changes over the *zero*-th (*I 0*) type within the mutually recursive family *Tree23Code* with values of type *MetaVar* in its holes.

*Computing Changes.* Computing a *Change* codes *MetaVar* from a source and a destination elements of type *Fix* codes *ix* follows exactly the structure as we saw previously in Section 2: extract the pair of contexts and fix unbound metavariables in a postprocessing step. We are still assuming an efficient oracle *buildOracle* *s d* :: *Oracle* codes, that determines whether or not an arbitrary *t* is a subtree of a fixed *s* and *d* indexed by *n*, where:

**type** *Oracle* codes =  $\forall j . \text{Fix codes } j \rightarrow \text{Maybe Int}$

*buildOracle* :: *Fix* codes *i*  $\rightarrow$  *Fix* codes *i*  $\rightarrow$  *Oracle* codes

The core of computing a change is in the extraction of the deletion and insertion contexts (*extract* function, Section 2). We must take care to avoid the problem we encountered in our previous implementation: a subtree that occurs in both the source and destination trees, but also occurs as the subtree of another common subtree (Figure 3) may result in unbound metavariables. We have shown how to fix this with a postprocessing step of the resulting change. That is still the case, but we now collect additional information from the context extraction before postprocessing.

Looking at the type of *Oracle*, we see it will only share recursive positions by construction. We use the *ForceI* type to bring this fact on the type level. That is, we are only sharing *recursive* positions so far. We also use the indexed product type (*.\**) to carry along information.

**data** (*.\**) *f g* *x* = *f x* *.\** *g x*

**data** *ForceI* :: (*Nat*  $\rightarrow$  *\**)  $\rightarrow$  *Atom*  $\rightarrow$  *\** **where**

*ForceI* :: *f i*  $\rightarrow$  *ForceI f* (*I i*)

Defining the generic *txExtract* function is simple. We check whether a given *x* is a subtree of the source and destination trees by consulting the oracle. If so, we return a hole with the subtree annotated; if *x* is not a common subtree we extract the topmost constructor and recurse over its recursive positions.



```

txPostprocess :: Tx codes (ForceI (Const Int *: Fix codes)) (I ix)
  → Tx codes (ForceI (Const Int *: Fix codes)) (I ix)
  → Change (ForceI (Const Int)) (I ix)
txPostprocess del ins =
  let okvars = varsOf del ∩ varsOf ins
  -- We have to txJoin the results since keepOrDrop returns a Tx
  in Change (txJoin (utxMap (keepOrDrop okvars) del))
    (txJoin (utxMap (keepOrDrop okvars) ins))
where
  -- traverses a Tx and puts all the variables in a Set.
  varsOf :: Tx codes (ForceI (Const Int *: Fix codes)) (I ix) → Set Int
  -- check whether a variable is in both contexts and decides
  keepOrDrop okvars (ForceI (Const mvar) *: subtree)
    | var ∈ okvars = TxHole (ForceI (Const mvar))
    | otherwise   = na2tx (NAI subtree)

```

Fig. 9. Post processing of contexts returning closed changes

```

txExtract :: Oracle codes
  → Fix codes ix
  → Tx codes (ForceI (Const Int *: Fix codes)) (I ix)
txExtract wcs x = case wcs x of
  Just i   → TxHole (ForceI (Const i *: x))
  Nothing  → let Tag c p = sop (unFix x)
             in TxPeel c (mapNP (elimNA TxOpq (txExtract wcs)) p)

```

Postprocessing works by traversing the result of the extracted contexts. In case we need to keep a given tree and forget that it was shared we convert it to a context with  $na2tx$ . Recall the reason for wanting to keep only the metavariables that occur in both the insertion and deletion contexts is to prevent any undefined variable when applying our patches, which would break correctness. More technically, the  $txPostprocess$  function groups the metavariables of each context in a set and computes the intersection of such sets, then maps over its arguments replacing the  $Const\ Int\ *: Fix\ codes$  hole by either  $Const\ Int$ , if the  $Int$  belongs in the set, or by translating the  $Fix\ codes$  with  $na2tx \circ NA_I$ . The implementation is shown in Figure 9.

At this point, given two trees  $a$  and  $b$  of type  $Fix\ codes\ ix$ , we can extract and post-process both the deletion and insertion contexts, of type  $Tx\ codes\ (ForceI\ (Const\ Int))\ (I\ ix)$ . These are just like a value of type  $Fix\ codes\ ix$  with holes of type  $Const\ Int$  exclusively in recursive positions. This is the generic version of the  $changeTree23$  function presented in Section 2:

```

change :: Fix codes ix → Fix codes ix → Change codes (ForceI (Const Int)) (I ix)
change x y = let wcs = buildOracle x y
             in txPostprocess (txExtract wcs x) (txExtract wcs y)

```

Recall that this function will only produce closed changes. That is, a deletion and a insertion context that have the same set of variables. Intuitively, this means that every variable that is declared is used and vice-versa.

*Minimizing the Changes: Computing Patches.* The next step is to minimize the changes coming from *change* function, yielding a *patch*. The generic counterpart of *Patch23* from Section 2 is given by:

**type** *Patch* codes at = *Tx* codes (*Change* codes *MetaVar*) at

Firstly, we have to identify the *spine*, that is, the constructors that are present in both the deletion and insertion contexts. This is done by splitting a big change into the *greatest common prefix* of the insertion and deletion contexts and the smaller changes inside. When going over two *TxPeel* we must check that the constructors are the same, i.e., have the same index. We use *testEquality* from *Data.Type.Equality* to check for type index equality and inform the compiler of that fact by matching on *Refl*.

```
txGCP :: Tx codes φ at → Tx codes ψ at → Tx codes (Tx codes φ ∗ Tx codes ψ) at
txGCP (TxOpq x) (TxOpq y)
  | x ≡ y      = TxOpq x
  | otherwise = TxHole (TxOpq x ∗ TxOpq y)
txGCP (TxPeel cx px) (TxPeel cy py)
  = case testEquality cx px of
      Nothing → TxHole (TxPeel cx px ∗ TxPeel cy py)
      Just Refl → TxPeel cx (mapNP (uncurry' txGCP) (zipNP px py))
txGCP a b = TxHole (a ∗ b)
```

The *txGCP* can be just like a generalized *zip* function, but instead of stopping whenever one of its arguments has been consumed and forgetting the other, it stores the rest of the other argument. It is *greatest* in the sense that it consumes as many constructors as possible and resorts to *TxHole* when it cannot progress further.

We know, from Section 2 that we cannot simply take the result of *change*, compute its *greatest common prefix* and be done with it. This would yield a patch with potentially malformed changes. The *txGCD* function is not aware of metavariables and might break their scoping.

Refining the result of *txGCP* is conceptually simple. All we have to do is bubble up the changes to a point where they are all *closed*. We start by creating some machinery to distinguish the open changes from the closed changes in the result of a *txGCP*. Then we define a traversal that shall look at those tags and decide whether more constructors should be pushed into the changes or not. This is done by the *closure* function.

Tagging open and closed changes is done with the indexed sum type. We use *InL* to mark the *open* changes and *InR* for the *closed* changes.

```
data Sum f g x = InL (f x) | InR (g x)
either' :: (f x → r x) → (g x → r x) → Sum f g x → r x
either' a b (InL fx) = a fx
either' a b (InR gx) = b gx
```

The *isClosed* predicate is responsible for deciding how to tag a change.

```
isClosed :: Change codes at → Sum (Change codes MetaVar) (Change codes MetaVar) at
isClosed (Change del ins)
  | variables ins ≡ variables del = InR (Change del ins)
  | otherwise                     = InL (Change del ins)
```

The *Sum*  $f$   $g$ , of kind *Atom*  $\rightarrow *$ , comes in handy for it can be passed as an argument to *Tx*. This enables us to map our predicate over an arbitrary patch  $p$ :

*txMap isClosed*  $p :: Tx$  codes (*Sum* (*Change* codes *MetaVar*) (*Change* codes *MetaVar*)) at

The final *closure* function is defined with an auxiliary function, *closure'*. This auxiliary function receives a patch with tagged changes and tries to eliminate all the *open changes*, tagged with *InL*. We do so by finding the smallest closed change that binds the required variables. If the *closure'* function cannot find the constructor that binds all variables, it tags the patch as an *open change* altogether. The first three cases are trivial:

*closure'*  $:: Tx$  codes (*Sum* (*Change* codes *MetaVar*) (*Change* codes *MetaVar*)) at  
 $\rightarrow Sum$  (*Change* codes *MetaVar*) (*Tx* codes (*Change* codes *MetaVar*)) at  
*closure'* (*TxOpq*  $x$ ) = *InR* (*TxOpq*  $x$ )  
*closure'* (*TxHole* (*InL*  $oc$ )) = *InL*  $oc$   
*closure'* (*TxHole* (*InR*  $cc$ )) = *InR*  $cc$

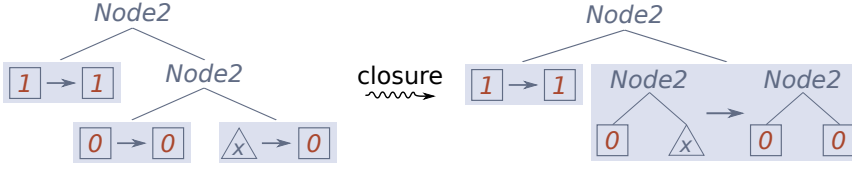
The interesting case of the *closure'* function is the *TxPeel* pattern, where we first try to compute the closures for the fields of the constructor and check whether all these fields contain only closed changes. If that is the case, we are done. If some fields contain open changes, however, the *mapNPM fromInR* fails with a *Nothing* and we must massage some data. This is akin to a simple distributive law for *Tx*, defined below.

*distr*  $:: Tx$  codes (*Change* codes  $\varphi$ ) at  $\rightarrow Change$  codes  $\varphi$  at  
*distr spine* = *Change* (*txJoin* (*txMap* *chgDel* *spine*))  
(*txJoin* (*txMap* *chgIns* *spine*))

The difference between *distr* and the *Nothing* clause in *closure'* is that in the later we are handling *NP* (*Tx* codes (*Change* codes  $\varphi$ )), i.e., a sequence of *Tx* instead of a single one. Consequently, we need some more code.

*closure'* (*TxPeel*  $cx$   $px$ )  
= **let**  $aux = mapNP$  *closure'*  $px$   
**in case** *mapNPM fromInR*  $aux$  **of**  
*Just*  $np \rightarrow InR$  (*TxPeel*  $cx$   $np$ ) -- everything is closed.  
-- some changes are open. Try pushing  $cx$  down the changes in  $px$  and  
-- see if this closes them, then.  
*Nothing*  $\rightarrow$  **let**  $chgs = mapNP$  (*either'* *InL* (*InR*  $\circ$  *distr*))  $aux$   
 $dels = mapNP$  (*either'* *chgDel* *chgDel*)  $chgs$   
 $inss = mapNP$  (*either'* *chgIns* *chgIns*)  $chgs$   
 $tmp = Change$  (*TxPeel*  $cx$   $dels$ ) (*TxPeel*  $cx$   $inss$ )  
**in if** *isClosed*  $tmp$   
**then** *InR* (*TxHole*  $tmp$ )  
**else** *InL* (*Change*  $tmp$ )

In the code above, *aux* is a sequence of either open changes or patches. The local *dels* and *inss* are defined as the a sequence deletion and insertion contexts from *aux*, regardless if they come from open or closed changes. This allows us to assemble a new, larger, change (*tmp*). Finally, we check whether this larger change is closed or not. We recall the illustration in Figure 6, repeated below, for a graphical intuition.



It is worth mentioning that in some corner cases, *closure'* will *forget* about certain contractions, producing a patch that has a larger domain than the original. This is arguable a good thing, as we want patches that can be applied to more things.

To finish it up, we wrap *closure'* within a larger function that always returns a *Tx* with all changes being *closed*. The necessary observation is that if we pass a given *tx* to *closure'* such that *distr tx* is closed, then *closure'* will always return a *InR* value. In our case, the *txPostprocess* is in place precisely to provided that guarantee, hence, the *error* is unreachable.

```
closure :: Tx codes (Sum (Change codes) (Change codes)) at
  → Patch codes at
closure = either' (const $ error "no closure exists") id
```

The final *diff* function is then assembled by using the closure of the greatest common prefix of the change the comes from the *change* function. In order to further enlarge the domain of our patches we add a small additional step where we replace the opaque values in the spine with copies.

```
diff :: Fix codes ix → Fix codes ix → Patch codes (I ix)
diff x y = let Change del ins = change x y
  in closure $ txRefine TxHole mkCpy
    $ txMap isClosed
    $ txGCP del ins
```

The *txRefine* simply traverses the *Tx* and refines the holes and opaques into *Txs* using the provided functions. In our case we leave the holes unchanged and replace the occurrences of *TxOpq* by a new *copy* change.

```
txRefine :: (∀ at . f at → Tx codes g at)
  → (∀ k . Opq k → Tx codes g (K k))
  → Tx codes f at → Tx codes g at
```

**Applying Patches.** Patch application follows closely the scheme sketched in for 2-3-trees (Section 2). There is one main difference, though. Our changes are now placed in the leaves of our spine and can be applied locally.

Our final *applyChange* will be responsible for receiving a change and a tree and instantiate the metavariables by matching the tree against the deletion context then substituting this valuation in the insertion context. Its type is given by:

```
applyChange :: Change codes MetaVar at
  → NA (Fix codes) at
  → Maybe (NA (Fix codes) at)
```

We are now left to match the spine with a value of *NA (Fix codes)*, and leave the changes paired up with the corresponding local elements they must be applied to. This is similar to the *txGCP*, and can be implemented by it. We must extract the greatest common prefix of the spine and the *Tx*

that comes from translating the  $NA$  (*Fix codes*) value but must make sure that the leaves have all *TxHoles* on the left.

```
txZipEl :: Tx codes  $\varphi$  at  $\rightarrow NA$  (Fix codes) at  $\rightarrow Maybe$  (Tx codes ( $\varphi$  ::  $NA$  (Fix codes)))
txZipEl tx el = txMapM (uncurry' checkIsHole) $ txGCP tx (tx2na el)
```

**where**

```
checkIsHole :: Tx codes  $\varphi$  at  $\rightarrow Tx$  codes  $\psi$  at  $\rightarrow Maybe$  ( $\varphi$  ::  $NA$  (Fix codes) at)
checkIsHole (TxHole  $\varphi$ ) el = ( $\varphi$ ::) <$> na2tx el
checkIsHole _ _ = Nothing
```

Finally, we define our application function. First we check whether the spine matches the element. If that is the case, we apply the changes, which are already paired with the parts of the element they must be applied to:

```
apply :: Patch codes ix  $\rightarrow Fix$  codes ix  $\rightarrow Maybe$  (Fix codes ix)
apply patch el = txZipEl patch el  $\gg$  return  $\circ$  txMapM (uncurry' applyChange)
```

Whenever a patch  $p$  can be applied to an element  $x$ , that is,  $apply\ p\ x$  returns *Just*  $y$  for some  $y$ , we say that  $p$  is *applicable* to  $x$ .

### 3.3 Defining the Generic Oracle

We conclude the generic algorithm with the implementation of the generic oracle which answers whether a tree is a common subtree of the source and destination of a patch. In this section we take the example from Section 2.4 and generalize the construction.

The first step is annotating our trees with cryptographic hashes [Menezes A. J. and Vanstone [n. d.]]. Essentially transforming our trees into *merkle trees* [Merkle 1988]. This technique is more commonly seen in the security and authentication context [Miller et al. 2014; Miraldo et al. 2018], and is similar in spirit to *hash-consing* [Filliâtre and Conchon 2006]. Luckily, the generic programming machinery that is already at our disposal enables us to create *merkle trees* generically quite easily. The generics-mrsop provide some attribute grammar [Knuth 1990] functionality, in particular the computation of synthesized attributes arising from a fold. The *synthesize* function is just like a catamorphism, but we decorate the tree with the intermediate results at each node, rather than only using them to compute the final outcome. This enables us to decorate each node of a *Fix codes* with a unique identifier (as shown in Figure 7, for example) by running the generic *decorate* function, defined below.

```
newtype AnnFix x codes i = AnnFix (x i, Rep (AnnFix x codes) (Lkup codes i))
decorate :: Fix codes i  $\rightarrow AnnFix$  (Const Digest) codes i
decorate = synthesize authAlgebra
merkleRoot :: AnnFix (Const Digest) codes i  $\rightarrow Digest$ 
merkleRoot (AnnFix (Const r, _)) = r
```

Here, *AnnFix* is the cofree comonad, used to add a label to each recursive branch of our generic trees. In our case, this label will be the cryptographic hash of the concatenation of its subtree's hashes. We can easily fetch these hashes with the *merkleRoot* function. Figure 7 shows an example of an input and corresponding output of the *decorate* function, producing a *merkelized Tree23*. The *synthesize* generic combinator annotates each node of the tree with the result of the catamorphism called at that point with the given algebra. Our algebra is sketched in pseudo-Haskell below:

```

authAlgebra :: Rep (Const Digest) sum → Const Digest iy
authAlgebra rep = case sop rep of
  Tag c [p1, ..., pn] → Const ∘ sha256Concat
    $ [encode c, encode (getSNat @iy), p1, ..., pn]

```

We must append the index of the type in question, in this case `getSNat @iy`, to our hash computation to differentiate constructors of different types in the family represented by the same number. In the real implementation we must also pass around a constraint stating that constructors can be encoded and opaque values can be hashed.

Following the guidelines in Section 2.4, we now must traverse our merkelized tree and store all the digests we see in a Trie [Brass 2008].

```
mkSharingTrie :: AnnFix (Const Digest) codes j → Trie Word MetaVar
```

It is of paramount importance to avoid recomputing the merkle root of a tree  $x$  each time we wish to know whether it is a common subtree. Otherwise, we still end up with an exponential algorithm. The solution is quite simple: we use `AnnFix (Const Digest) codes` in the `txExtract` function and the type of our oracle, where `Fix codes` was used before. This provides access to the merkle root in constant time. After this small modification to our `Oracle`, allowing it to receive trees annotated with hashes we proceed to define the efficient `buildOracle` function.

```

type Oracle codes = ∀ j . AnnFix (Const Digest) codes j → Maybe Int
buildOracle :: Fix codes i → Fix codes i → Oracle codes
buildOracle s d = let s' = decorate s
                  d' = decorate d
                  in lookup (mkSharingTrie s' ∩ mkSharingTrie d') ∘ merkleRoot

```

It is true that the `buildOracle`, as defined above, could give us false positives caused by hash collisions. It is worth mentioning that when using a cryptographic hash, the chance of collision is negligible [Menezes A. J. and Vanstone [n. d.]] and we chose to ignore it. Regardless, it is not difficult to work around the possibility for collisions if we wanted. We could compute an intermediate `Trie` from `Word` to `Exists (Fix codes)` in the `mkSharingTrie` function and every time we find a potential collision we check the trees for structural equality. If equality check fails, a hash collision is found and the entry would be removed from the map. The `Exist` datatype simply encapsulate the `ix` type index from a `Fix codes ix` into an existential.

## 4 MERGING PATCHES

One of the main motivations for generic structure-aware diffing is being able to merge patches in a more structured fashion than using `diff3`, which considers changes to every line. In the past, structural merging has proven to be a difficult task [Miraldo et al. 2017; Vassena 2016] even for the easiest cases. This is due to the sub-optimal representations for structured patches. This section shows how our new structure for representing changes enables us to write a simple merge algorithm, offering both acceptable performance and a improvement over `diff3`. We will sketch the implementation of our algorithm here and evaluate its performance in Section 5.

The merging problem, illustrated in Figure 10, is the problem of computing a new patch,  $q // p$ , given two patches  $p$  and  $q$ . It consists in a patch that contains the changes of  $q$  adapted to work on a value that has already been modified by  $p$ . This is sometimes called the *transport* of  $q$  over  $p$  or the *residual* [Huet 1994] of  $p$  and  $q$ .

There is a class of patches that are trivial to merge: those that modify separate locations of a tree. If  $p$  and  $q$  are disjoint, then  $p // q$  can return  $p$  without further adaptations. Our algorithm

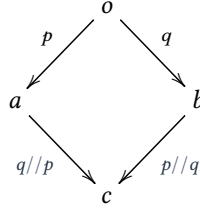


Fig. 10. Merge square

shall merge only disjoint patches, marking all other situations as a conflict. We choose to represent conflicts as a pair of overlapping patches.

**type** *Conflict* codes = *Patch* codes *∗*: *Patch* codes

**type** *PatchConf* codes = *Tx* codes (Sum (*Conflict* codes) (*Change* codes *MetaVar*))

In practice, it may be desirable to record further meta-information to facilitate conflict resolution.

Our merging operator, (*//*), receives two patches and returns a patch possibly annotated with conflicts. We do so by matching the spines, and carefully inspecting any changes where the spines differ.

(*//*) :: *Patch* codes at → *Patch* codes at → *PatchConf* codes at

The intuition here is that  $p // q$  must preserve the intersection of the spines of  $p$  and  $q$  and reconcile the differences whenever one of the patches has a change. Note that it is impossible to have completely disjoint spines since  $p$  and  $q$  are applicable to at least one common element. Using the *greatest common prefix* function defined previously, we can zip together the shared spines, pushing the resolution down to the leaves:

$p // q = \text{txMap } (\text{uncurry}' \text{ reconcile}) \$ \text{txGCP } p \ q$

Here, the *reconcile* function shall check whether the disagreeing parts are disjoint, i.e., either one of the two changes is the identity or they perform the exactly same change. If that is the case, the *reconcile* function returns its first argument. In fact, this is very much in line with the properties of a residual operator [Huet 1994].

*reconcile* :: *Patch* codes at → *Patch* codes at → Sum (*Conflict* codes) (*Change* codes) at

*reconcile*  $p \ q$

|  |  |
|--|--|
| <i>patchIden</i> $p \vee \text{patchIden } q = \text{InR } \$ \text{distr } p$ |  |
| <i>patchEquiv</i> $p \ q$  | = <i>InR</i> $\$ \text{copy}$            |
| <i>otherwise</i>   | = <i>InL</i> $\$ \text{Conflict } p \ q$ |

We see the code for *reconcile* closely follows the definition of disjointness above—one of the patches must be the identity or they are equal. The *patchIden* functions checks whether all changes in that patch are copies and *patchEquiv* checks if two patches are  $\alpha$ -equivalent. Taking a closer look at the *reconcile* function, we see it follows the three identity laws from residual theory. The first branch agrees from the two identity laws from residual theory that state that  $p // \text{id} \equiv p$  and  $\text{id} // p \equiv \text{id}$ , whereas the second branch follows the third identity law, which states that  $p // p \equiv \text{id}$ , meaning that applying a patch over something that has been modified by this very patch amounts to not changing anything.



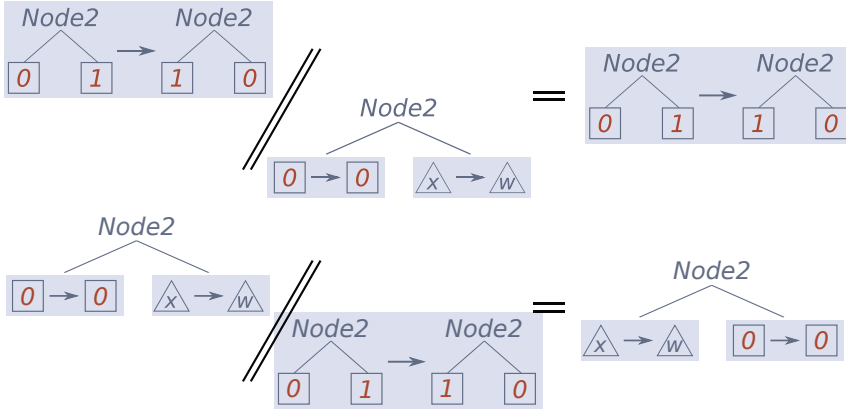


Fig. 11. Example of a merge square where the first residual is obtained by not changing the patch and the second is computed by applying a patch to another patch, transporting the changes.

Our trivial merge algorithm returns a conflict for non-disjoint patches, but this does not mean that it is impossible to merge them in general. Although a full discussion is out of the scope of this paper, there are a number of non-disjoint patches that can still be merged. These non-trivial merges can be divided in two main situations: (A) there is no action needed even though patches are not disjoint, and (B) the relevant parts of a patch can be transported to operate on different parts tree automatically. In Figure 11 we illustrate situations (A) and (B) in the merge square for two non-disjoint patches. In the top subfigure we see the residual returning the patch unaltered (case A). In this example, the patch in the ‘nominator’ position is a simple swap of subtrees. This swap operation can be applied to every possible result of applying the second patch in the ‘denominator’ of the residual. As a result, it computing the residual is easy: we simply return the ‘nominator’ patch. In the second subfigure, however, the situation is reversed. The ‘denominator’ patch must be applied to the ‘nominator’ – yielding a new patch that has the expected behavior. In future work, we hope to identify the precise conditions under which two non-disjoint patches can be merged in this way.

## 5 EXPERIMENTS

We have conducted two experiments over a number of Lua [Jerusalimschy et al. 1996] source files. We obtained these files data by mining the top Lua repositories on GitHub and extracting all the merge conflicts recorded in their history. Next, we ran two experiments over this data: a performance test and a merging test. We chose the Lua programming language for two reasons. First, there is a Haskell parser for Lua readily available on Hackage and, secondly, due to a performance bug [GHC Trac 2018] in GHC we are not able to instantiate our generic algorithm to more complex abstract syntax trees, such as that of C.

*Performance Evaluation.* In order to evaluate the performance of our implementation<sup>1</sup> we have timed the computation of the two diffs,  $\text{diff } o \ a$  and  $\text{diff } o \ b$ , for each merge conflict  $a, o, b$  in our dataset. In order to ensure that the numbers we obtained are valid and representative of a real execution trace we timed the execution time of parsing the files and running  $\text{length } o \ \text{encode } o \ \text{uncurry } \text{diff}$  over the parsed files, where *encode* comes from *Data.Serialize*. Besides ensuring that

<sup>1</sup><https://github.com/VictorCMiraldo/hs-digems>

| Repository         | Commits | Contributors | Total Conflicts | Trivial Conflicts |
|--------------------|---------|--------------|-----------------|-------------------|
| awesome            | 9289    | 250          | 5               | 0                 |
| busted             | 936     | 39           | 9               | 0                 |
| CorsixTH           | 3207    | 64           | 25              | 8                 |
| hawkthorne-journey | 5538    | 61           | 158             | 27                |
| kong               | 4669    | 142          | 163             | 11                |
| koreader           | 6418    | 89           | 14              | 2                 |
| luakit             | 4160    | 62           | 28              | 2                 |
| luarocks           | 2131    | 51           | 45              | 3                 |
| luvit              | 2862    | 68           | 4               | 1                 |
| nn                 | 1839    | 177          | 3               | 0                 |
| Penlight           | 730     | 46           | 6               | 3                 |
| rnn                | 622     | 42           | 6               | 1                 |
| snabb              | 8075    | 63           | 94              | 6                 |
| tarantool          | 12299   | 82           | 33              | 2                 |
| telegram-bot       | 729     | 50           | 5               | 0                 |
| total              |         |              | 598             | 66                |

Fig. 12. Lua repositories mined from *GitHub*

the patch is fully evaluated, the serialization also mimics what would happen in a real version control system since the patch would have to be saved to disk. After timing approximately 1200 executions from real examples we have plotted the data over the total number of constructors for each source-destination pair. In Figure 13 we see two plots: on the left we have plotted 70% of our dataset in more detail whereas on the right we show the full plot. The results were expected given that we seen how  $\text{diff } x \ y$  runs in  $O(n + m)$  where  $n$  and  $m$  are the number of constructors in  $x$  and  $y$  abstract syntax trees, respectively. Confirming our analysis with empirical results further strengthens our algorithm as a practical implementation of structured differencing, even though it stands about one order or magnitude slower than purely textual differencing tools. One could combine both approaches and only fallback to the more expensive structural merging tool when the textual merging fails. In Figure 13 we see that around 90% of our dataset falls within a one second runtime.

*Merging Evaluation.* We have also performed a preliminary evaluation of the simple merging algorithm presented in Section 4. After collecting all the merge commits from the GitHub repositories, we selected those that `git merge` failed to solve automatically, which we call a conflict, and attempted to use our structured merge instead. When this merge succeeded, we checked that the resulting merge square (Figure 10) commutes as expected. In this way, we were able to solve a total of 66 conflicts automatically, amounting to 11% of all the conflicts we encountered. We consider these initial numbers to be encouraging: even a naive merge algorithm on structured changes manages to outperform the current state of the art. We expect that a more refined notion of merging may improve these results further.

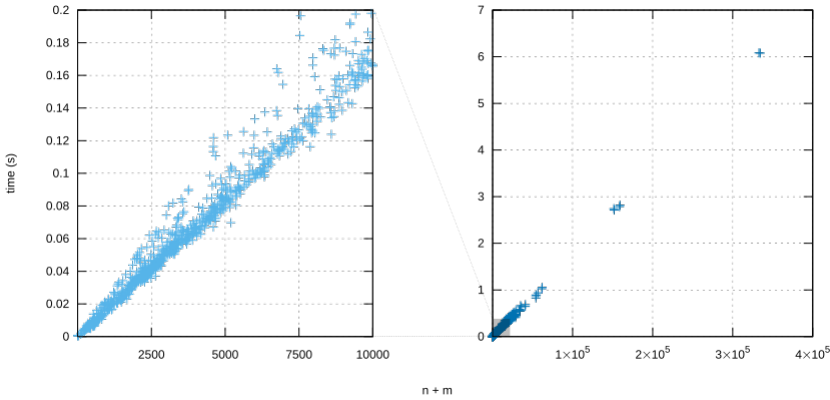


Fig. 13. Plot of the time for diffing two lua files over the total AST nodes

### 5.1 Threats to Validity

There are two main threats to the validity of our empirical results. Firstly, we are diffing and merging *abstract* syntax trees, hence ignoring comments and formatting. There would be no extra effort in handling these issues, beyond recording them explicitly in a more concrete syntax tree and adapting our parser to produce such trees. Nevertheless, one might expect a slightly lower success rate since we are ignoring formatting changes altogether. Secondly, a significant number of developers prefer to rebase their branches instead of merging them. Therefore, we may have missed a number of important merge conflicts that are no longer recorded, as rebasing erases history. Our merge algorithm might be able to resolve some of these conflicts automatically—but there is no way to establish this. Nevertheless, we could still have gathered the *Pull Requests* merged using rebase or squash on GitHub. This is because GitHub keeps a hidden branch that refers to the state before the merge available.

## 6 DISCUSSION AND CONCLUSION

The results from Section 5 are very encouraging. We see that our diffing algorithm has competitive performance and our trivial merging operation is capable of merging changes where `git merge` fails. Yet there is still plenty of work to be done.

### 6.1 Future Work

*Controlling Sharing.* One interesting direction for further work is how to control the sharing of subtrees. As it stands, the differencing algorithm will share every subtree that occurs in both the source and destination files. This can lead to undesirable behavior. For example, we may not want to share *all* occurrences of a variable within a program, but rather only share occurrences of a variable with the same binder. That is, sharing should respect the scope variables. A similar question arises with constants – should all occurrences of the number `1` be shared?

There are a variety of options to customize the sharing behavior of our algorithm. One way to do so would allow the definition of a custom oracle that is scope-aware. By hashing both the identifier name and its binder, we can ensure that variables are not shared over scope boundaries. Another option would be to consider abstract syntax trees that make the binding structure of variables explicit.

*Better Merge Algorithm.* The merging algorithm presented in Section 4 only handles trivial cases. Being able to merge patches that are not disjoint is the subject of ongoing research. The problem seems related to unification, residual systems, and rewriting systems. We hope that relating the merging problem to these settings might help nail down the necessary conditions for merging to succeed. One would expect that it would have some resemblance to a pushout, as in pointed out by Mimram and Di Giusto [Mimram and Giusto 2013]. This would come hand-in-hand with more meta-theoretical work of formalizing our representation of patches and our algorithms in dependently typed language. Another interesting addendum to a better merging algorithm is the ability to define domain specific strategies to solve conflicts.

*Extending the Generic Universe.* Our prototype is built on top of `generics-mrsop`, a generic programming library for handling mutually recursive families in the sums of products style. With recent advances in generic programming [Serrano and Miraldo 2018], we might be able to extend our algorithm to handle mutually recursive families that have GADTs.

## 6.2 Related Work

The hashing techniques used in this paper is reminiscent of *hash-consing* [Filliâtre and Conchon 2006], a technique to share values that are structurally equal. Hash-consing is done by maintaining a global hash-table during run-time, which keeps track of the values that have already been created and can be shared. There are two main differences to our situation, however: (A) we cannot use hash tables for it would use too much space, and, (B) we must detect which subtrees are shared within two fixed trees instead of sharing the values in memory, which is the main objective of *hash-consing*. On another hand, the problem of minimizing finite acyclic deterministic automata has efficient solutions in the literature [Bubenzer 2014], which could be seen as a different way of defining the *which common subtree* oracle.

On the diffing side, related work can be classified in the treatment of types. The untyped tree differencing problem was introduced in 1979 [Tai 1979] as a generalization of the longest common subsequence problem [Bergroth et al. 2000]. There has been a significant body of work on the untyped tree differencing problem [Akutsu et al. 2010; Demaine et al. 2007; Klein 1998], but these results do not transport to the typed setting: the transformations that are computed are not guaranteed to produce well-typed trees.

The first datatype generic algorithm was presented by Lempink and Löh [Lempink et al. 2009], which was later extended by Vassena [Vassena 2016]. Their work consists largely in using the same algorithm as `diff` on the flattened representation of a tree. The main observation is that basic operations (insertion, deletion and copy) can be shown to be well-typed when operating on these flattened representations. Although one could compute differences with reasonably fast algorithms, merging these changes is fairly difficult and in some cases might be impossible [Vassena 2016]. Miraldo et al. [Miraldo et al. 2017] take a slightly different approach, defining operations that work directly on tree shaped data. Using this approach, changes become easier to merge but harder to compute. Both bodies of work follow the same general idea as the untyped variants: compute all possible patches and select the ‘best’ patch from these alternatives. As we have already mentioned (Section 1), this is not an optimal strategy. For one, the number of patches can grows exponentially if one is not careful to design a dynamic or memoized program. Secondly, defining the *best* patch using insertions, deletions and copies is impossible without further heuristics.

The work of Asenov et al. [Asenov et al. 2017] is also untyped, but uses a different technique for finding the diff: it flattens trees and embellishes the resulting lists with additional annotations, and then uses the UNIX `diff` tool to compute patches. Finally, it transports the changes back to the tree-shaped datatypes using the annotations that were added. The authors identify a number

of interesting situations that occur when merging tree differences. The gumtree [Falleri et al. 2014] project, explores a similar line of work, but uses its own algorithm for computing graph transformations between untyped representations of abstract syntax trees.

There have been several different approaches to formalizing a theory of patches. The version control system darcs [Roundy 2005] was one of the first to present a more formal theory of patches, but the patches themselves were still line-based. Mimram and De Giusto [Mimram and Giusto 2013] have developed a theoretical model of line-based patches in a categorical fashion. This has inspired the version control system pi.jul. Swierstra and Löh [Swierstra and Löh 2014] have proposed using separation logic to define a meta-theory of patches and merging. Finally, Angiuli et al. [Angiuli et al. 2014] describe a patch theory based on homotopy type theory.

Shifting to the language-theoretic point of view, two lines of similar research must be mentioned here. Firstly, incremental parsing [Wagner and Graham 1998] must also represent a change within a tree. It is different from structured differencing for incremental parsing does not have a notion of *change*. The change always happen at the point the user is editing the file, and its representation is trivial. The hard challenge is in applying parsing rules to partial input. Secondly, the work on Grammar-Based Tree Compression [Lohrey 2015] could be seen as a variant of the differencing problem. There are significant differences, however. In GBTC one must maintain one copy of every subtree, by definition. In our case, we want to abstract away the common subtrees and identify copy opportunities without storing the subtrees involved.

### 6.3 Conclusions

Throughout this paper we have developed an efficient type-directed algorithm for computing structured differences for a large class of algebraic datatypes, namely, mutually recursive families. This class of types can represent the abstract syntax tree of most programming languages and, hence, our algorithm can be readily instantiated to compute the difference between programs written in these languages. We have validated our implementation by computing diffs between Lua [Jerusalimschy et al. 1996] source files obtained from various repositories on GitHub; the algorithm’s run-time is competitive, and even a naive merging algorithm already offers a substantial improvement over existing technology, for the former is tree-based and the latter is line-based. Together, these results demonstrate both a promising direction for further research and a novel application of the generic programming technology that is readily available in today’s functional languages.

### REFERENCES

- Tatsuya Akutsu, Daiji Fukagawa, and Atsuhiko Takasu. 2010. Approximating Tree Edit Distance through String Edit Distance. *Algorithmica* 57, 2 (2010), 325–348. <https://doi.org/10.1007/s00453-008-9213-z>
- Carlo Angiuli, Edward Morehouse, Daniel R. Licata, and Robert Harper. 2014. Homotopical Patch Theory. In *Proceedings of the 19th ACM SIGPLAN International Conference on Functional Programming (ICFP ’14)*. ACM, New York, NY, USA, 243–256. <https://doi.org/10.1145/2628136.2628158>
- Dimitar Asenov, Balz Guenat, Peter Müller, and Martin Otth. 2017. Precise Version Control of Trees with Line-Based Version Control Systems. In *Proceedings of the 20th International Conference on Fundamental Approaches to Software Engineering - Volume 10202*. Springer-Verlag New York, Inc., New York, NY, USA, 152–169. [https://doi.org/10.1007/978-3-662-54494-5\\_9](https://doi.org/10.1007/978-3-662-54494-5_9)
- L. Bergroth, H. Hakonen, and T. Raita. 2000. A survey of longest common subsequence algorithms. In *String Processing and Information Retrieval, 2000. SPIRE 2000. Proceedings. Seventh International Symposium on*. 39–48.
- Philip Bille. 2005. A survey on tree edit distance and related problems. *Theor. Comput. Sci* 337 (2005), 217–239.
- Peter Brass. 2008. *Advanced Data Structures* (1 ed.). Cambridge University Press, New York, NY, USA.
- Johannes Bubenzer. 2014. Cycle-aware minimization of acyclic deterministic finite-state automata. *Discrete Applied Mathematics* 163 (2014), 238 – 246. <https://doi.org/10.1016/j.dam.2013.08.003> Stringology Algorithms.
- Edsko de Vries and Andres Löh. 2014. True Sums of Products. In *Proceedings of the 10th ACM SIGPLAN Workshop on Generic Programming (WGP ’14)*. ACM, New York, NY, USA, 83–94. <https://doi.org/10.1145/2633628.2633634>

- Erik D. Demaine, Shay Mozes, Benjamin Rossman, and Oren Weimann. 2007. An Optimal Decomposition Algorithm for Tree Edit Distance. In *Proceedings of the 34th International Colloquium on Automata, Languages and Programming (ICALP 2007)*. Wroclaw, Poland, 146–157.
- Richard A. Eisenberg and Stephanie Weirich. 2012. Dependently Typed Programming with Singletons. In *Proceedings of the 2012 Haskell Symposium (Haskell '12)*. ACM, New York, NY, USA, 117–130. <https://doi.org/10.1145/2364506.2364522>
- Jean-Rémy Falleri, Floréal Morandat, Xavier Blanc, Matias Martinez, and Martin Monperrus. 2014. Fine-grained and accurate source code differencing. In *ACM/IEEE International Conference on Automated Software Engineering, ASE '14, Vasteras, Sweden - September 15 - 19, 2014*. 313–324. <https://doi.org/10.1145/2642937.2642982>
- Jean-Christophe Filliâtre and Sylvain Conchon. 2006. Type-safe Modular Hash-consing. In *Proceedings of the 2006 Workshop on ML (ML '06)*. ACM, New York, NY, USA, 12–19. <https://doi.org/10.1145/1159876.1159880>
- GHC Trac. 2018. Memory usage exploding for complex pattern matching. (2018). <https://ghc.haskell.org/trac/ghc/ticket/14987#no2>.
- G rard Huet. 1994. Residual theory in  $\lambda$ -calculus: a formal development. *Journal of Functional Programming* 4, 3 (1994), 371–394. <https://doi.org/10.1017/S0956796800001106>
- J. W. Hunt and M. D. McIlroy. 1976. *An Algorithm for Differential File Comparison*. Technical Report CSTR 41. Bell Laboratories, Murray Hill, NJ.
- Roberto Ierusalimschy, Luiz Henrique de Figueiredo, and Waldemar Celes Filho. 1996. Lua: An Extensible Extension Language. *Software: Practice and Experience* 26, 6 (1996), 635–652. [https://doi.org/10.1002/\(SICI\)1097-024X\(199606\)26:6<635::AID-SPE26>3.0.CO;2-P](https://doi.org/10.1002/(SICI)1097-024X(199606)26:6<635::AID-SPE26>3.0.CO;2-P)
- Philip N. Klein. 1998. Computing the Edit-Distance Between Unrooted Ordered Trees. In *Proceedings of the 6th Annual European Symposium on Algorithms (ESA '98)*. Springer-Verlag, London, UK, UK, 91–102.
- Donald E. Knuth. 1990. The Genesis of Attribute Grammars. In *Proceedings of the International Conference WAGA on Attribute Grammars and Their Applications*. Springer-Verlag, London, UK, UK, 1–12. <http://dl.acm.org/citation.cfm?id=645938.671208>
- Elco Lempersink, Sean Leather, and Andres L h. 2009. Type-safe Diff for Families of Datatypes. In *Proceedings of the 2009 ACM SIGPLAN Workshop on Generic Programming (WGP '09)*. ACM, New York, NY, USA, 61–72.
- Markus Lohrey. 2015. Grammar-Based Tree Compression. In *Developments in Language Theory*, Igor Potapov (Ed.). Springer International Publishing, Cham, 46–57.
- Paul van Oorschot Menezes A. J. and Scott A. Vanstone. [n. d.]. *Handbook of Applied Cryptography* (boca raton, xiii, 780, 1997 ed.). CRC Press.
- Ralph C. Merkle. 1988. A Digital Signature Based on a Conventional Encryption Function. In *Advances in Cryptology – CRYPTO '87*, Carl Pomerance (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 369–378.
- Andrew Miller, Michael Hicks, Jonathan Katz, and Elaine Shi. 2014. Authenticated Data Structures, Generically. In *Proceedings of the 41st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '14)*. ACM, New York, NY, USA, 411–423. <https://doi.org/10.1145/2535838.2535851>
- Samuel Mimram and Cinzia Di Giusto. 2013. A Categorical Theory of Patches. *CoRR* abs/1311.3903 (2013). arXiv:1311.3903 <http://arxiv.org/abs/1311.3903>
- Victor Cacciari Miraldo, Harold Carr, Alex Kogan, Mark Moir, and Maurice Herlihy. 2018. Authenticated Modular Maps in Haskell. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Type-Driven Development (TyDe 2018)*. ACM, New York, NY, USA, 1–13. <https://doi.org/10.1145/3240719.3241790>
- Victor Cacciari Miraldo, Pierre- variste Dagand, and Wouter Swierstra. 2017. Type-directed diffing of structured data. In *Proceedings of the 2nd ACM SIGPLAN International Workshop on Type-Driven Development*. ACM, 2–15.
- Victor Cacciari Miraldo and Alejandro Serrano. 2018. Sums of products for mutually recursive datatypes: the appropriation-ist s view on generic programming. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Type-Driven Development*. ACM, 65–77.
- David Roundy. 2005. Darcs: Distributed Version Management in Haskell. In *Proceedings of the 2005 ACM SIGPLAN Workshop on Haskell (Haskell '05)*. ACM, New York, NY, USA, 1–4. <https://doi.org/10.1145/1088348.1088349>
- Alejandro Serrano and Victor Cacciari Miraldo. 2018. Generic Programming of All Kinds. In *Proceedings of the 11th ACM SIGPLAN International Symposium on Haskell (Haskell 2018)*. ACM, New York, NY, USA, 41–54. <https://doi.org/10.1145/3242744.3242745>
- Wouter Swierstra and Andres L h. 2014. The Semantics of Version Control. In *Proceedings of the 2014 ACM International Symposium on New Ideas, New Paradigms, and Reflections on Programming & Software (Onward! '14)*. 43–54.
- Kuo-Chung Tai. 1979. The Tree-to-Tree Correction Problem. *J. ACM* 26, 3 (July 1979), 422–433. <https://doi.org/10.1145/322139.322143>
- Marco Vassena. 2016. Generic Diff3 for Algebraic Datatypes. In *Proceedings of the 1st International Workshop on Type-Driven Development (TyDe 2016)*. ACM, New York, NY, USA, 62–71.

- Tim A. Wagner and Susan L. Graham. 1998. Efficient and Flexible Incremental Parsing. *ACM Trans. Program. Lang. Syst.* 20, 5 (Sept. 1998), 980–1013. <https://doi.org/10.1145/293677.293678>
- Brent A. Yorgey, Stephanie Weirich, Julien Cretin, Simon Peyton Jones, Dimitrios Vytiniotis, and José Pedro Magalhães. 2012. Giving Haskell a Promotion. In *Proceedings of the 8th ACM SIGPLAN Workshop on Types in Language Design and Implementation (TLDI '12)*. ACM, New York, NY, USA, 53–66. <https://doi.org/10.1145/2103786.2103795>