

Type-Safe Generic Differencing of Mutually Recursive Families

Victor Cacciari Miraldo

Type-Safe Generic Differencing of Mutually Recursive Families

Getypeerde Generieke Differentiatie van Wederzijds
Rekursieve Datatypes
(met een samenvatting in het Nederlands)

Proefschrift

ter verkrijging van de graad van doctor aan de Universiteit Utrecht
op gezag van de rector magnificus, prof. dr. H.R.B.M. Kummeling,
ingevolge het besluit van het college voor promoties in het openbaar
te verdedigen op maandag 5 oktober 2020 des ochtends te 11:00 uur

door

Victor Cacciari Miraldo

geboren op 16 oktober 1991
te São Paulo, Brazilië

Promotor: Prof.dr. G.K. Keller

Copromotor: Dr. W.S. Swierstra

Dit proefschrift werd (mede) mogelijk gemaakt met financiële steun van de Nederlandse Organisatie voor Wetenschappelijk Onderzoek (NWO), project Revision Control of Structured Data (612.001.401).

To my Mother, Father and Brother

ABSTRACT

The UNIX `diff` tool – which computes the differences between two files in terms of a set of copied lines – is widely used in software version control. The fixed *lines-of-code* granularity, however, is sometimes too coarse and obscures simple changes, i.e., renaming a single parameter triggers the whole line to be seen as *changed*. This may lead to unnecessary conflicts when unrelated changes occur on the same line. Consequently, it is difficult to merge such changes automatically.

In this thesis we discuss two novel approaches to structural differencing, generically – which work over a large class of datatypes. The first approach defines a type-indexed representation of patches and provides a clear merging algorithm, but it is computationally expensive to produce patches with this approach. The second approach addresses the efficiency problem by choosing an extensional representation for patches. This enables us to represent transformations involving insertions, deletions, duplication, contractions and permutations which are computable in linear time. With the added expressivity, however, comes added complexity. Consequently, the merging algorithm is more intricate and the patches can be harder to reason about.

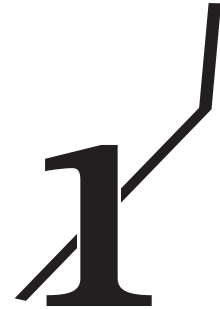
Both of our approaches can be instantiated to mutually recursive datatypes and, consequently, can be used to compare elements of most programming languages. Writing the software that does so, however, comes with additional challenges. To address this we have developed two new libraries for generic programming in Haskell.

Finally, we empirically evaluate our algorithms by a number of experiments over real conflicts gathered from `GitHub`. Our evaluation reveals that at least 26% of the conflicts that developers face on a day-to-day basis could have been automatically merged. This suggests there is a benefit in using structural differencing tools as the basis for software version control.

CONTENTS

ABSTRACT	iii
1 INTRODUCTION	1
1.1 Contributions and Outline	5
2 BACKGROUND	7
2.1 Differencing and Edit Distance	8
2.1.1 String Edit Distance and UNIX <code>diff</code>	9
2.1.2 Classic Tree Edit Distance	12
2.1.3 Shortcomings of Edit-Script Based Approaches	15
2.1.4 Synchronizing Changes	16
2.1.5 Literature Review	18
2.2 Generic Programming	20
2.2.1 GHC Generics	20
2.2.2 Explicit Sums of Products	22
2.2.3 Discussion	24
3 GENERIC PROGRAMMING WITH MUTUALLY RECURSIVE TYPES	27
3.1 The <code>generics-mrsop</code> library	28
3.1.1 Explicit Fixpoints with Codes	29
3.1.2 Mutual Recursion	32
3.1.3 Practical Features	40
3.1.4 Example: Well-Typed Classical Tree Differencing	43
3.2 The <code>generics-simplistic</code> Library	46
3.2.1 The Simplistic View	46
3.2.2 Mutual Recursion	48
3.2.3 The (Co)Free (Co)Monad	52
3.2.4 Practical Features	54
3.3 Discussion	57
4 STRUCTURAL PATCHES	59
4.1 The Type of Patches	61
4.1.1 Functorial Patches	61
4.1.2 Recursive Changes	65
4.2 Merging Patches	68
4.3 Computing <code>Patch_{ST}</code>	72
4.3.1 Naive enumeration	72
4.3.2 Translating from <code>gdiff</code>	74
4.4 Discussion	77

5	PATTERN-EXPRESSION PATCHES	79
5.1	Changes	80
5.1.1	A Concrete Example	80
5.1.2	Representing Changes Generically	88
5.1.3	Meta Theory	90
5.1.4	Computing Changes	95
5.2	The Type of Patches	100
5.2.1	Computing Closures	104
5.2.2	The <i>diff</i> Function	106
5.2.3	Aligning Changes	107
5.2.4	Summary	115
5.3	Merging Aligned Patches	116
5.4	Discussion and Further Work	129
6	EXPERIMENTS	133
6.1	Data Collection	134
6.2	Performance	134
6.3	Synchronization	137
6.3.1	Threats to Validity	143
6.4	Discussion	143
7	DISCUSSION	145
7.1	The Future of Structural Differencing	145
7.2	Future Work	147
7.3	Concluding Remarks	148
A	SOURCE-CODE AND DATASET	149
A.1	Source-Code	149
A.2	Dataset	149
	SAMENVATTING	163
	CURRICULUM VITAE	167
	ACKNOWLEDGEMENTS	169



INTRODUCTION

Version Control is essential for most distributed collaborative work. It enables contributors to operate independently and later combine their work. For that, though, it must address the situation where two developers changed a piece of information in different ways. One option is to lock further edits until a human decides how to reconcile the changes, regardless of the changes. Yet, many changes can be reconciled *automatically*.

Software engineers usually rely on version control systems to help with this distributed workflow. These tools keep track of the changes performed to the objects under version control, computing changes between old and new versions of an object. When time comes to reconcile changes, it runs a *merge* algorithm that decides whether the changes can be synchronized or not. At the heart of this process is (A) the representation of changes, usually denoted a *patch*, (B) the computation of a *patch* between two objects and (C) the ability to detect whether two patches are in conflict.

Maintaining software as complex as an operating system with as many as several thousands contributors is a technical feat made possible thanks, in part, to a venerable Unix utility: `UNIX diff` [46]. It computes the line-by-line difference between two textual files, determining the smallest set of insertions and deletions of lines to transform one file into the other. In other words, it tries to share as many lines between source and destination as possible. This is, in fact, the most widespread representation for *patches*, used by tools such as Git, Mercurial and Darcs.

The limited grammar of changes used by the `UNIX diff` works particularly well for programming languages that organize a program into lines of code. For example, consider the modification in Figure 1.1, where the insertions do not interfere with the rest of the program.

```

    sum := 0;
+   prod := 1;
    for (i in is) {
        sum += i;
+   prod *= i;
    }

```

FIGURE 1.1: *Modification that extends an existing for-loop to not only compute the sum of the elements of an array, but also compute their product.*

However, the bias towards *lines* of code may lead to (unnecessary) conflicts when considering other programming languages. For instance, consider the following diff between two Haskell functions that adds a new argument to an existing function:

```

- head []          = error "?!"
- head (x :: xs) = x
+ head []          d = d
+ head (x :: xs) d = x

```

This modest change impacts all the lines of the function's definition, even though it affects relatively few elements of the abstract-syntax.

The line-based bias of the diff algorithm may lead to unnecessary *conflicts* when considering changes made by multiple developers. Consider the following innocuous improvement of the original `head` function, which improves the error message raised when the list is empty:

```

head []          = error "Expecting a non-empty list."
head (x :: xs) = x

```

Trying to apply the patch above to this modified version of the `head` function will fail, as the lines do not match – even if both changes modify distinct parts of the declaration in the case of non-empty lists.

The inability to identify more fine grained changes in the objects being compared is a consequence of the *by line* granularity of *patches*. Ideally, however, the objects under comparison should dictate the granularity of change to be considered. This is precisely the goal of *structural differencing* tools.

If we reconsider the example above, we could give a more detailed description of the modification made to the `head` function by describing the changes made to the constituent declarations and expressions:

```

head []          {+d+} = error {"?!"} {"Expect..."}
head (x :: xs) {+d+} = x

```

There is more structure here than mere lines of text. In particular, the granularity is at the abstract-syntax level. It is worthwhile to note that this problem also occurs in languages which tend to be organized in a line-by-line manner. Modern languages which contain any degree of object-orientation will also group several abstract-syntax elements on the same line. Take the Java function below,

```
public void test(obj) {
    assert(obj.size(), equalTo(5));
}
```

Now consider a situation where one developer updated the test to require the size of `obj` to be 6, but another developer changed the function that makes the comparison, resulting in the two orthogonal versions below;

```
public void test(obj) {                public void test(obj) {
    assert(obj).hasSize(5);              assert(obj.size(), equalTo(6));
}
```

It is arguable that the desired *synchronized* version can incorporate both changes, calling `assert(obj).hasSize(6)`. Combining these changes would be impossible without access to information about the old and new state of *individual abstract-syntax elements*. Simple line-based information is insufficient, even in line-oriented languages.

Differencing and synchronization algorithms tend to follow a common framework – compute the difference between two values of some type *a*, and represent these changes in some type, *Patch a*. The *diff* function *computes* the differences between two values of type *a*, whereas *apply* attempts to transform a value according to the information stored in the *Patch* provided to it.

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

A definition of *Patch a* which has access to information about the structure of *a* enables us to represent changes at a more refined granularity. In Chapters 4 and 5 we discuss two different definitions of *Patch*, both capturing changes at the granularity of abstract-syntax elements.

Note that the *apply* function is inherently partial, for example, when attempting to delete data which is not present applying the patch will fail. 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 [7, 31] using `xml` or `json` to represent abstract syntax trees, where the result of applying a patch may produce ill-typed results, i.e., schema violations.

UNIX *diff* [46] follows this very framework too, but for the specific type of lines of text, taking *a* to be `[String]` and *Patch a* to be a series of insertions, deletions and copies of lines. A naive implementation would produce patches by enumerating all possibilities

that transform the source into the destination and then chooses the best such patch. The UNIX `diff` computes its patches in a more complex and efficient manner, but follows the above method as its *specification*. There have been several attempts at generalizing these results to handle arbitrary datatypes [114, 27, 87, 55], including our own attempt discussed in Chapter 4. All of these follow the same recipe: enumerate all combinations of insertions, deletions and copies that transform the source into the destination and choose the ‘best’ one. Consequently, they also suffer from the same drawbacks as classic edit-distance – which include non-uniqueness of the best solution and slow algorithms. We will discuss them in more detail in Section 2.1.1.

Once we have the *diff* and *apply* functions at hand, we move on to the *merge* function, which is responsible for synchronizing two different changes into a single one, when they are compatible. Naturally, we can only merge patches that alter *disjoint* parts of the AST. Hence, the merge function must be partial, returning a conflict whenever patches change the same part of the tree in different ways.

merge :: *Patch* *a* → *Patch* *a* → *Either Conflicts (Patch a)*

A realistic merge function should naturally distribute conflicts to their specific locations inside the merged patch and still try to synchronize non-conflicting parts of the changes. This is orthogonal to our objective, however. The abstract idea is still the same: two patches can either be reconciled fully or there exist conflicts between them.

The success rate of the *merge* function – that is, how often it is able to reconcile changes – can never be 100%. There will always be changes that require human intervention to be synchronized. Nevertheless, the quality of the synchronization algorithm directly depends on the expressivity of the *Patch* datatype. If *Patch* provides information solely on which lines of the source have changed, there is little we can merge. Hence, we want values of type *Patch* to carry information about the structure of *a*. Naturally though, we do not want to build domain specific tools for each programming language for which we wish to have source files under version control – which would be at least impractical. A better option is to use a *generic representation*, which can be used to encode arbitrary programming languages, and describe the *Patch* datatype generically.

Structural differencing is a good example of the need for generic programming. Here, we would like to have our differencing algorithms working over arbitrary abstract syntax trees while maintaining the type-safety that a language like Haskell provides: we encode these ASTs as algebraic datatypes and write our differencing algorithm to operate over these algebraic datatypes. This added safety means that all the manipulations we perform on the patches are guaranteed to never produce ill-formed elements, which is a clear advantage over using something like XML to represent our data, even though there exist differencing tools that use XML as their underlying representation for data. We refer to these as *untyped* tree differencing algorithms in contrast to the *typed* approach, which guarantees type safety by construction.

The Haskell type-system is expressive enough to enable one to write *typed* generic algorithms. These algorithms, however, can only be applied to datatypes that belong to the set of types handled by the generic programming library of choice. For example, the `regular` [82] approach is capable of handling types which have a *regular* recursive structure – lists, n -ary trees, etc. –, but cannot represent nested types, for example. In Section 2.2 we will give an overview of existing approaches to generic programming in Haskell. No library, however, was capable of handling mutually recursive types – which is the universe of datatypes that context free languages belong in – in a satisfactory manner. This means that to explore differencing algorithms for various programming languages we would have to first develop the generic programming functionality necessary for it. Happily, Haskell’s type system has evolved enough since the initial efforts on generic programming for mutually recursive types (`multirec` [112]), enabling us to write significantly better libraries, as we will discuss in Chapter 3.

1.1 CONTRIBUTIONS AND OUTLINE

This thesis documents a number of peer-reviewed contributions, namely:

- a) Chapter 3 discusses the `generics-mrsop` [75] library, which offers combinator-based generic programming for mutually recursive families. This work came out of close collaboration with Alejandro Serrano on a variety of generic programming topics.
- b) Chapter 4 is derived from a paper [74] published with with Pierre-Évariste Dagand. We worked closely together to define a type-indexed datatype used to represent changes in a more structured way than edit-scripts. Chapter 4 goes further into developing a merging algorithm and exploring different ways to compute patches given two concrete values. The code we present in Chapter 4 is loosely based on Van Putten’s translation of our Agda repository to Haskell as part of his Master thesis work [92].
- c) Chapter 5 is the refinement of our paper [76] on an efficient algorithm for computing patches, where we tackle the problems from Chapter 4 with a different representation for patches altogether.

Other contributions that have not been peer-reviewed include:

- d) Chapter 3 discusses the `generics-simplistic` library, a different approach to generic programming that overcomes an important space leak in the Haskell compiler, which rendered `generics-mrsop` unusable in large, real-world, examples.

- e) Chapter 5 introduces a merging algorithm and Chapter 6 discusses its empirical evaluation over a dataset of real conflicts extracted from `GitHub`.



BACKGROUND

The most popular tool for computing differences between two files is UNIX `diff` [46], which works by comparing files in a *line-by-line* basis and attempts to match lines from the source file to lines in the destination file. For example, consider the two files below:

1	<code>res := 0;</code>	1	<code>print("summing up");</code>
2	<code>for (i in is) {</code>	2	<code>sum := 0;</code>
3	<code>res += i;</code>	3	<code>for (i in is) {</code>
4	<code>}</code>	4	<code>sum += i;</code>
		5	<code>}</code>

Lines 2 and 4 in the source file, on the left, match lines 3 and 5 in the destination. These are identified as copies. The rest of the lines, with no matches, are marked as deletions or insertions. In this example, lines 1 and 3 in the source are deleted and lines 1,2 and 4 in the destination are inserted.

This information about which lines have been *copied*, *deleted* or *inserted* is then packaged into an *edit-script*: a list of operations that transforms the source file into the destination file. For the example above, the edit-script would read something like: delete the first line; insert two new lines; copy a line; delete a line; insert a line and finally copy the last line. The output we would see from UNIX `diff` would show deletions prefixed with a minus sign and insertions prefixed with a plus sign. Copies have no prefix. In our case, it would look something like:

```
- res := 0;
+ print("summing up");
+ sum := 0;
  for (i in is) {
-   res += i;
+   sum += i;
  }
```


The edit-scripts produced by the UNIX `diff` contain information about transforming the source into the destination file by operating exclusively at the *lines-of-code* level. Computing and representing differences in a finer granularity than *lines-of-code* is usually done by parsing the data into a tree and later flattening said tree into a list of nodes, where one then reuses existing techniques for computing differences over lists, i.e., think of printing each constructor of the tree into its own line. This is precisely how most of the classic work on tree edit distance computes tree differences (Section 2.1.2).

Recycling linear edit distance into tree edit distance, however, comes with its drawbacks. Linear differencing uses *edit-scripts* to represent the differences between two objects. Edit-scripts are composed of atomic operations, which traditionally include at least *insert*, *delete* and *copy*. These scripts are later interpreted by the application function, which gives the semantics to these operations. The notion of *edit distance* between two objects is defined as the minimum possible cost associated with an *edit-script* between them, where cost is some metric which is often context dependent. One major drawback, for example, is the least cost edit-script is chosen arbitrarily in some situations, namely, when it is not unique. This makes the results computed by these algorithms hard to predict, and consequently, so is the result of merging patches.

The algorithms computing edit-scripts must either return an approximation of the least cost edit-script or check countless ambiguous choices to return the optimal one. Finally, manipulating edit-scripts in an untyped fashion, say, for instance in order to merge them, might produce ill-typed trees – as in *not abiding by a schema* – as a result [107]. We can get around this last issue by writing edit-scripts in a typed form [55], but this requires some non-trivial generic programming techniques to scale.

The first half of this chapter introduces some of the classical *edit-script* based algorithms whereas the second half of presents the state-of-the-art of the generic programming ecosystem in Haskell.

2.1 DIFFERENCING AND EDIT DISTANCE

The *edit distance* between two objects is defined as the least possible cost of an edit-script that transforms the source object into the target object – in its simplest form, it can be seen as the cost of the edit-script with the least insertions and deletions. Computing edit-scripts is often referred to as *differencing* objects. Where edit distance computation is only concerned with how *similar* one object is to another, *differencing*, on the other hand, is actually concerned with how to transform one objects into another. Although very closely related, these do make up different problems. In the biology domain [2, 42, 67], for example, one is concerned solely in finding similar structures in a large set of structures, whereas in software version control systems manipulating and combining differences is important.

The wide applicability of differencing and edit distances leads to a variety of cost notions, edit-script operations and algorithms for computing them [16, 14, 84]. In this section we will review some of the important notions and background work on edit distance. We start by looking at the string edit distance (Section 2.1.1) and then generalize this to untyped trees (Section 2.1.2), as it is classically portrayed in the literature, which is reviewed in Section 2.1.5.

2.1.1 STRING EDIT DISTANCE AND UNIX `DIFF`

In this section we look at two popular notions of edit distance. The *Levenshtein Distance* [56, 14], for example, works well for detecting spelling mistakes [80] or measuring how similar two languages are [103]. It considers insertions, deletions and substitutions of characters as its edit operations. The *Longest Common Subsequence (LCS)* [14], on the other hand, considers insertions, deletions and copies as edit operations and is better suited for identifying shared sequences between strings.

LEVENSHTEIN DISTANCE The Levenshtein distance regards insertions, deletions and substitutions of characters as edit operations, which can be modeled in Haskell by the *EditOp* datatype below. Each of those operations has a predefined *cost* metric.

```
data EditOp = Ins Char | Del Char | Subst Char Char
cost :: EditOp → Int
cost (Ins _)      = 1
cost (Del _)      = 1
cost (Subst c d) = if c == d then 0 else 1
```

These individual operations are then grouped into a list, usually denoted an *edit-script*. The *apply* function, below, gives edit-scripts a denotational semantics by mapping them to partial functions over *Strings*.

```
apply :: [EditOp] → String → Maybe String
apply [] [] = Just []
apply (Ins c : ops) ss = (c :) <$> apply ops ss
apply (Del c : ops) (s : ss) = guard (c == s) >> apply ops ss
apply (Subst c d : ops) (s : ss) = guard (c == s) >> (d :) <$> apply ops ss
apply _ _ = Nothing
```

The *cost* metric associated with these edit operations is defined to force substitutions to cost less than insertions and deletions. This ensures that the algorithm looking for the list of edit operations with the minimum cost will prefer substitutions over deletions and insertions.

```

lev :: String → String → [EditOp]
lev [] [] = []
lev (x : xs) [] = Del x : lev xs []
lev [] (y : ys) = Ins y : lev [] ys
lev (x : xs) (y : ys) = let i = Ins y : lev (x : xs) ys
                        d = Del x : lev xs (y : ys)
                        s = Subst x y : lev xs ys
                        in minimumBy cost [i, d, s]

```

FIGURE 2.1: Definition of the function that returns the edit-script with the minimum Levenshtein Distance between two strings.

We can compute the *edit-script*, i.e. a list of edit operations, with the minimum cost quite easily with a brute-force and inefficient specification, illustrated in Figure 2.1.

```

levenshteinDist :: String → String → Int
levenshteinDist s d = cost (head (lev s d))

```

Note that although the Levenshtein distance is unique, the edit-script witnessing it is *not*. Consider the case of lev “ab” “ba” for instance. All of the edit-scripts below have cost 2, which is the minimum possible cost.

```

lev “ab” “ba” ∈ [ [ Del 'a' , Subst 'b' 'b' , Ins 'a' ]
                  , [ Ins 'b' , Subst 'a' 'a' , Del 'b' ]
                  , [ Subst 'a' 'b' , Subst 'b' 'a' ] ]

```

From an edit distance point of view, this is not an issue. The Levenshtein distance between “ab” and “ba” is 2, regardless of the edit-script. But from an operational point of view, i.e., transforming one string into another, this ambiguity poses a problem. The lack of criteria to favor one edit-script over another means that the result of the differencing algorithm is hard to predict. Consequently, developing a predictable diff and merging algorithm becomes a difficult task.

LONGEST COMMON SUBSEQUENCE

Given our context of source-code version-control, we are rather interested in the *Longest Common Subsequence (LCS)*, which is a restriction of the Levenshtein distance and forms the specification of the UNIX `diff` [46] utility.

If we take the *lev* function and modify it in such a way that it only considers identity substitutions, that is, *Subst* *x* *y* with *x* ≡ *y*, we end up with a function that computes

```

lcs :: [String] → [String] → [EditOp]
lcs [] [] = []
lcs (x : xs) [] = Del x : lcs xs []
lcs [] (y : ys) = Ins y : lcs [] ys
lcs (x : xs) (y : ys) = let i = Ins y : lcs (x : xs) ys
                        d = Del x : lcs xs (y : ys)
                        s = if x == y then [Cpy x : lcs xs ys] else []
                        in minimumBy cost (s + [i, d])

```

FIGURE 2.2: Specification of the UNIX *diff*.

the classic longest common subsequence. Note that this is different from the longest common substring problem, as subsequences need not be contiguous.

UNIX *diff* [46] computes a solution to the LCS problem between two *files*, seen as a list of *strings*, opposed to a list of *characters*. Hence, the edit operations become:

```

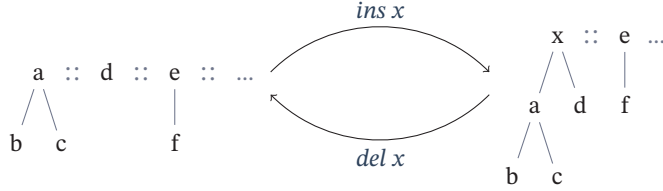
data EditOp = Ins String | Del String | Cpy String
cost :: EditOp → Int
cost (Ins _) = 1
cost (Del _) = 1
cost (Cpy _) = 0

```

The application function is analogous to the *apply* for the Levenshtein distance. The computation of the minimum cost edit-script, however, is not. We must ensure to issue a *Cpy* only when both elements are the same, as illustrated in Figure 2.2.

Running the *lcs x y* function, Figure 2.2, will yield an *edit-script* that enables us to read out one longest common subsequence of *x* and *y*. Note that the ambiguity problem is still present, however to a lesser degree than with the Levenshtein distance. For example, there are only two edit-scripts with minimum cost on *lcs* [“a”, “b”] [“b”, “a”]. This, in fact, is a general problem with any *edit-script* based approach.

The original UNIX *diff* implementation was based on Hirschberg’s dynamic algorithm [44], which uses a memoized *lcs* to avoid recomputing sub-problems and has a quadratic runtime. The current implementation is based on Myers algorithm [79] and runs in $\mathcal{O}(d(n+m))$, where *n* and *m* are the size of the input files and *d* is the edit distance between them. Actual implementations also employ a number of algorithmic tricks to make it more performant, for instance, it is common to hash the data being compared to have amortized constant time comparison. There is also a number of heuristics that tend to perform well in practice. One example is the *diff --patience* algorithm [24], that will emphasize the matching of lines that appear only once in the source and destination files.

FIGURE 2.3: Insertion and Deletion of node x , with arity 2 on a forest.

2.1.2 CLASSIC TREE EDIT DISTANCE

Tree edit-distance is a generalization of the (linear) edit-distance problem. Instead of computing a distance between two lists of values, we are interested in a distance between two *trees* of values. The classical algorithms [3, 27, 50, 16, 10, 21] consider *untyped* trees – directed acyclic graphs where each vertex has at most one parent – as the objects under scrutiny. We call them *untyped* in the sense that they do not abide by any schema: nodes can have a label and an arbitrary number of children, opposed to a *typed* tree which must abide by a given schema, i.e., it can be seen as a value of a family of ADTs in Haskell, where the type signatures provide the schema.

There is an added degree of freedom that comes from considering trees instead of lists, and this carries over to the choice of edit operations. Suddenly, there are more edit operations one could use to create edit-scripts. To name a few, we can have flattening insertions and deletions, where the children of the deleted node are inserted or removed in-place in the parent node, or node relabeling. This degree of variation is responsible for the high number of different approaches and techniques we see in practice [32, 41, 31, 84, 34], as addressed in Section 2.1.5.

Basic tree edit distance [27], however, considers only node insertions, deletions and copies. The cost function is borrowed entirely from string edit distance together with the longest common subsequence function, that instead of working with $[a]$ will now work with $[Tree]$. Figure 2.3 illustrates insertions and deletions of (untyped) labels on a forest. The interpretation of these edit operations as actions on forests is shown in Figure 2.4.

We label these approaches as *untyped* because there exist edit-scripts that yield non-well formed trees. For example, imagine l is a label with arity 2 – supposed to receive two arguments. Now consider the edit-script $Ins\ l : []$, which will yield the tree `Node l []` once applied to the empty forest. If the objects under differencing are required to abide by a certain schema, such as abstract syntax trees for example, this becomes an issue. Granted we could define *apply* to take arities into account, this is not the case for the classical algorithms in the literature. This issue becomes particularly relevant when one needs to manipulate patches independently of the objects they have been created from.

```

data EOp = Ins Label | Del Label | Cpy Label
data Tree = Node Label [ Tree ]
arity :: Label → Int
apply :: [EOp] → [Tree] → Maybe [Tree]
apply [] [] = Just []
apply (Cpy l : ops) ts = apply (Ins l : Del l : ops) ts
apply (Del l : ops) (Node l' xs : ts) = guard (l ≡ l') >> apply ops (xs ++ ts)
apply (Ins l : ops) ts = (λ(args, rs) → Node l args : rs) ∘ takeDrop (arity l)
                                <$> apply ops ts
apply _ _ = Nothing

```

FIGURE 2.4: Definition of *apply* for tree edit operations.

Imagine a merge function that needs to construct a patch based on two other patches. A wrong implementation of said merge function can yield invalid trees for some given schema. In the context of abstract-syntax, this could be unparseable programs.

It is possible to use the Haskell type system to our advantage and write *EOp* in such a way that it is guaranteed to return well-typed results. Labels will be the different constructors of the family of types in question and their arity comes from how many fields each constructor expects. Edit-scripts will then be indexed by two lists of types: the types of the trees it consumes and the types of the trees it produces. We will come back to this in more detail in Section 3.1.4, where we review the approach of Lempink and Löh [55] at adapting this untyped framework to be type-safe by construction.

Although edit-scripts (Figure 2.4) provide a very intuitive notion of local transformations over a tree, there are many different edit-scripts that perform the same transformation: the order of insertions and deletions does not matter. This makes it hard to develop algorithms based solely on edit-scripts. The notion of *tree mapping* often comes in handy. It works as a *normal form* version of edit-scripts and represents only the nodes that are either relabeled or copied. We must impose a series of restrictions on these mappings to maintain the ability to produce edit-scripts out of it. Figure 2.5 illustrates four invalid and one valid such mappings.

Definition 2.1.1 (Tree Mapping). Let t and u be two trees, a tree mapping between t and u is an order preserving partial bijection between the nodes of a flattened representation of t and u according to their preorder traversal. Moreover, it preserves the ancestral order of nodes. That is, given two subtrees x and y in the domain of the mapping m , then x is an ancestor of y if and only if $m\ x$ is an ancestor of $m\ y$. We say that x is an ancestor of y if x is reachable from y proceeding exclusively from child to parent.

The tree mapping determines the nodes where either a copy or substitution must be performed. Everything else must be deleted or inserted and the order of deletions and

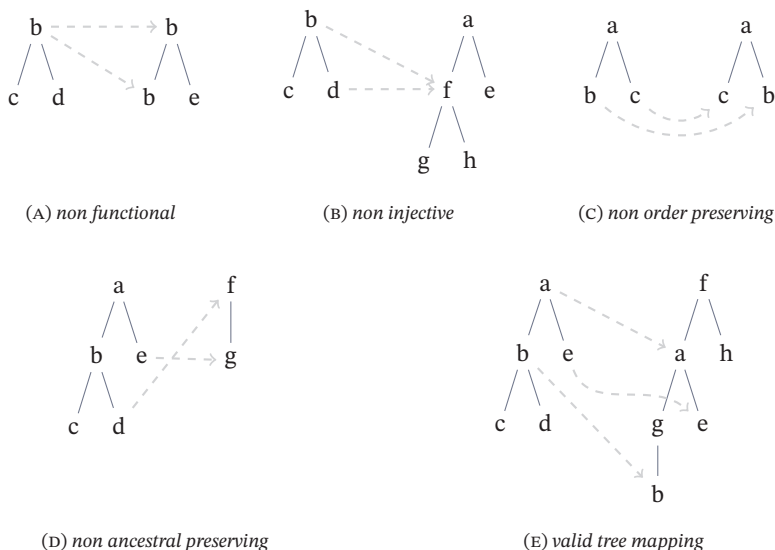


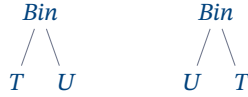
FIGURE 2.5: A number of invalid tree mappings with one valid example.

insertions is irrelevant, which removes the redundancy of edit-scripts. Nevertheless, the definition of tree mapping is still very restrictive: the “bijective mapping” does not enable trees to be duplicated or contracted, as seen in Figures 2.5(A) and 2.5(B); the “order preserving” does not enable trees to be permuted or moved across ancestor boundaries, as seen in Figures 2.5(C) and 2.5(D). These restrictions are there to ensure that one can always compute an edit-script from a tree mapping.

Most tree differencing algorithms start by producing a tree mapping and then extracting an edit-script from this. There are a plethora of design decisions on how to produce a mapping and often the domain of application of the tool will enable one to impose extra restrictions to attempt to squeeze maximum performance out of the algorithm. The `LaDiff` [22] tool, for example, works for hierarchically structured trees – used primarily for \LaTeX source files – and uses a variant of the LCS to compute matchings of elements appearing in the same order, starting at the leaves of the document. Tools such as `XyDiff` [23], used to identify changes in XML documents, use hashes to produce matchings efficiently.

2.1.3 SHORTCOMINGS OF EDIT-SCRIPT BASED APPROACHES

We argue that regardless of the process by which an edit-script is obtained, edit-scripts have inherent shortcomings when they are used to compare tree structured data. The first and most striking is that the use of heuristics to compute optimal solutions is unavoidable. Consider the tree-edit-scripts between the following two trees:



From an *edit distance* point of view, their distance is 2. This fact can be witnessed by two (propositionally) different edit-scripts: both $[Cpy\ Bin, Del\ T, Cpy\ U, Ins\ T]$ and $[Cpy\ Bin, Ins\ U, Cpy\ T, Del\ U]$ transform the target into the destination correctly. Yet, from a *differencing* point of view, these two edit-scripts are distinct. Do we care more about U or T ? What if U and T are also trees, but happen to have the same size (so that inserting one or the other yields edit-scripts with equal costs)? Ultimately, differencing algorithms that support no *swap* operation must choose to copy T or U arbitrarily. This decision is often guided by heuristics, which makes the result of different algorithms hard to predict and reason about. Moreover, the existence of this type of choice point inherently slows algorithms down since the algorithm *must decide* which tree to copy.

Another issue when dealing with edit-script is that they are type unsafe. It is quite easy to write an edit-script that produces an *ill-formed* tree, regardless of the schema. Even when writing the edit operations in a type-safe way [55] the synchronization of said changes is not guaranteed to be type-safe [107].

Finally, we must mention the lack of expressivity that comes from edit-scripts, from the *differencing* point of view. Consider the trees below,



Optimal edit-scripts oblige us to choose between copying A as the left or the right subtree. There is no possibility to represent duplications, permutations or contractions of subtrees. This means that a number of common changes, such as refactorings, yield edit-scripts with a very high cost even though a good part of the information being deleted or inserted should really have been copied. Even though there exists other edit-distances that support more edit operations, they are not very useful when adapted to trees. Take the Damerau-Leveshtein distance [25], which allows for the transposition of adjacent characters in a string, when instantiated to trees it would only allow for the transposition of labels that are adjacent in the preorder traversal of the tree.

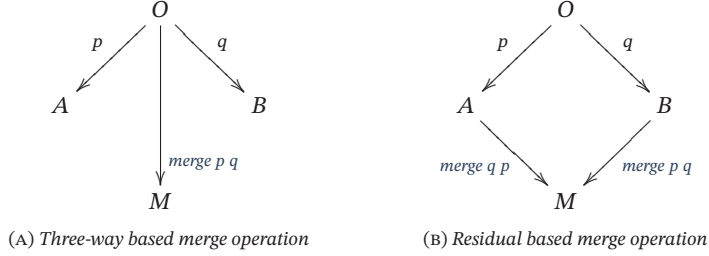


FIGURE 2.6: Two different ways to look at the merge problem.

2.1.4 SYNCHRONIZING CHANGES

When managing local copies of replicated data such as in software version control systems, one is inevitably faced with the problem of *synchronizing* [12] or *merging* changes – when an offline machine goes online with new versions, when two changes happened simultaneously, etc. The *synchronizer* is responsible for identify what has changed and reconcile these changes when possible. Most modern synchronizers operate over the diverging replicas and last common version, without knowledge of the history of the last common version – these are often denoted *state-based* synchronizers, as opposed to *operation-based* synchronizers, which access the whole history of modifications.

The `diff3` [101] tool, for example, is the most widely used synchronizer for textual data. It is a *state-based* synchronizer that calls the UNIX `diff` to compute the differences between the common ancestor and each diverging replica, then tries to produce an edit-script that when applied to the common ancestor produces a new file, containing the union of changes introduced in each individual replica. The algorithm itself has been studied formally [48] and there are proposals to extend it to tree-shaped data [57, 107].

Generally speaking, synchronization of changes p and q can be modeled in one of two ways. Either we produce one change that works on the common ancestor of p and q , as in Figure 2.6(A), or we produce two changes that act directly on the images of p and q , Figure 2.6(B). We often call the former a *three-way merge* and the later a *residual merge*.

Residual merges can pose a few technical challenges. For one, if we want the merge to form a (term rewriting) residual system [15] we must prove a number of non-trivial properties. Secondly, they tend to be harder to generalize to n -ary inputs. They do have the advantage of enabling one to model merges as pushouts [71], which could provide a desirable metatheoretical foundation on Category Theory.

Regardless of whether we choose a *three-way* or *residual* based approach, any state-based synchronizer will invariably have to deal with the problem of *aligning* the changes.

<pre> sum := 0; for (i in is) { sum := sum + i; } </pre>	<pre> res := 0; for (i in is) { res := res + i; } </pre>	<pre> res := 0; prod := 1; for (i in is) { res := res + i; prod := prod * i; } </pre>
(A) <i>Replica A</i>	(B) <i>Common ancestor, 0</i>	(C) <i>Replica B</i>

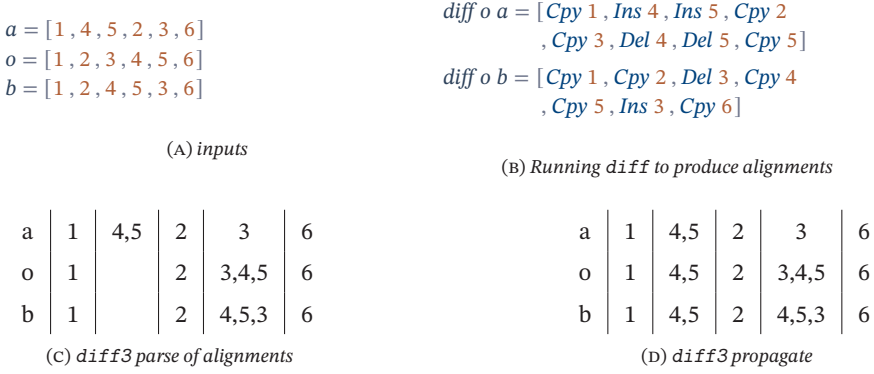
<pre> - res := 0; + sum := 0; for (i in is) { - res := res + i; + sum := sum + i; } </pre>	<pre> res := 0; + prod := 1; for (i in is) { res := res + i; + prod := prod * i; } </pre>
(D) <i>diff 0 A</i>	(E) <i>diff 0 B</i>

FIGURE 2.7: Two UNIX *diff* patches that diverge from a common ancestor.

That is, deciding which parts of the replicas are copies from the same piece of information in the common ancestor. For example, successfully synchronizing the replicas in Figure 2.7 depends in recognizing that the insertion of `prod := 1;` comes after modifying `res := 0;` to `sum := 0;`. This fact only becomes evident after we look at the result of calling the UNIX *diff* on each diverging replica – the copies in each patch identify which parts of the replicas are ‘the same’.

Figure 2.8 illustrates a run of *diff3* in a simple example, borrowed from Khanna et al. [48], where *a* swaps 2, 3 for 4, 5 in the original file but *b* moves 3 before 6. In a very simplified way, the first thing that happens if we run *diff3* in the inputs (Figure 2.8(A)) is that *diff3* will compute the longest common subsequences between the objects, essentially yielding the alignments it needs (Figure 2.8(B)). The next step is to put the copies side by side and understand which regions are *stable* or *unstable*. The stable regions are those where no replicas changed. In our case, it is on 1, 2 and 6 (Figure 2.8(C)). Finally, *diff3* can decide which changes to propagate and which changes are a conflict. In our case, the 4, 5 was only changed in one replica, so it is safe to propagate (Figure 2.8(D)).

Different synchronization algorithms will naturally offer slightly different properties, yet, one that seems to be central to synchronization is locality [48] – which is enjoyed by *diff3*. Locality states that well-separated changes of a given object can always be synchronized without conflicts. In fact, we argue this is the only property we can expect out of a general purpose generic synchronizer. The reason being that said synchronizer can rely solely on propositional equality of trees and structural disjointness as the criteria to establish changes as synchronizable. Any other criteria would require knowledge of the semantics of the data under synchronization. It is worth noting that although “well-separated changes” is difficult to define for an underlying list [48], tree

FIGURE 2.8: A simple *diff3* run.

shaped data has the advantage of possessing simpler such notions. We often refer to well-separated changes as *disjoint* changes.

2.1.5 LITERATURE REVIEW

With some basic knowledge of differencing and edit-distances under our belt, we briefly look over some of the relevant literature on the topic of tree differencing. Tai [102] was the first to consider edit distance between two trees, followed Zhang and Sasha [114]. More refined algorithms for computing the distance between two trees have been presented by Klein et al. [50] and Dulucq et al. [28]. Finally, Demaine et al. [27] presents an algorithm of cubic complexity and proves this is the best possible worst case. Zhang and Sasha's algorithm is still preferred in many practical scenarios, though. The more recent *RTED* [87] algorithm maintains the cubic worst case complexity and is comparable or faster than the other algorithms, rendering it the standard choice for computing tree edit distance based on the classic edit operations. In the case of unordered trees the best we can rely on are approximations [8, 9] since the problem is NP-hard [115].

Tree edit distance has seen multidisciplinary interest. From Computational Biology, where it is used to align phylogenetic trees and compare RNA secondary structures [2, 42, 67], all the way to intelligent tutoring systems where we must provide good hints to students' solutions to exercises by understanding how far they are from the correct solutions [85, 96]. In fact, from the *tree edit distance* point of view, we are only concerned with a number, the *distance* between objects, quantifying how similar they are.

From the perspective of *tree differencing*, on the other hand, we focus mainly on the edit operations and might want to perform computations such as composition and merg-

ing of differences. Naturally, however, the choice of edit operations heavily influences the complexity of the *diff* algorithm. Allowing a *move* operation already renders string differencing NP-complete [99]. Tree differencing algorithms, therefore, tend to run approximations of the best edit distance. Most of them still suffer from at least quadratic time complexity, which is prohibitive for most practical applications or are defined for domain specific data, such as the *latexdiff* [104] tool. A number of algorithms specific for XML and imposing different requirements on the schemas have been developed [88]. *LaDiff* [22], for example, imposes restrictions on the hierarchy between labels, it is implemented into the *DiffXML* [78] and *GumTree* [31] tools. *LaDiff* is responsible for computing an edit-script given tree matchings. A notable mention is the *XyDiff* [23], which uses hashes to compute matchings and, therefore, supports *move* operations maintaining almost linear complexity. This is the closest to our approach in Chapter 5. The *RWS-Diff* [34] uses approximate matchings by finding trees that are not necessarily equal but *similar*. This yields a robust algorithm, which is practical. Most of these techniques recycle list differencing and can be seen as some form of string differencing over the preorder (or postorder) traversal of trees, which has quadratic upper bound [40]. A careful encoding of the edit operations enables one to have edit-scripts that are guaranteed to preserve the schema of the data under manipulation [55].

When it comes to synchronization of changes [12], the algorithms are heavily dependent on the representation of objects and edit-scripts imposed by the underlying differencing algorithm. The *diff3* [101] tool, developed by Randy Smith in 1988, is still the most widely used synchronizer. It has received a formal treatment and specification [48] posterior to its development. Algorithms for synchronizing changes over tree shaped data include *3DM* [57] which merges changes over XML documents, *Harmony* [35], which works internally with unordered edge-labelled trees and is focused primarily on unordered containers and, finally, *FCDP* [54], which uses XML as its internal representation.

Also worth mentioning is the generalization of *diff3* to tree structured data using well-typed approaches due to Vassena [107], which supports that typed edit-scripts might not be the best underlying framework for this, as one needs to manually type-check the resulting edit-scripts.

Besides source-code differencing there is patch inference and generation tools. Some infer patches from human created data [49], whereas other, such as *Coccinelle* [6, 86], receive as input a number of diffs, P_0, \dots, P_n , that come from differencing many source and target files, $P_i = \text{diffs}_i t_i$. The objective then is to infer a common transformation that was applied everywhere. One can think of determining the *common denominator* of P_0, \dots, P_n . Refactoring and Rewriting Tools [68, 62] must also be mentioned. Some of these tools define each supported language AST separately [20, 51], whereas others [106] support a universal approach similar to *S-expressions*. They identify only parentheses, braces and brackets and hence, can be applied to a plethora of programming languages out-of-the-box.

2.2 GENERIC PROGRAMMING

We would like to consider richer datatypes than *lines-of-text*, without having to define separate *diff* functions for each of them. *Datatype-generic programming* provides a mechanism for writing functions by induction on the *structure* of algebraic datatypes [38]. A widely used example is the **deriving** mechanism in Haskell, which frees the programmer from writing repetitive functions such as equality [63]. A vast range of approaches are available as preprocessors, language extensions, or libraries for Haskell [95, 60].

The core idea behind generic programming is the fact that a number of datatypes can be described in a uniform fashion. Hence, if a programmer were to write programs that work over this uniform representation, these programs would immediately work over a variety of datatypes. In this section we look into two modern approaches to generic programming which are widely used, then discuss their design space and drawbacks.

2.2.1 GHC GENERICS

The `GHC.Generics` [59] library, which comes bundled with GHC since version 7.2, defines the representation of datatypes in terms of uniform *pattern functors*. Consider the following datatype of binary trees with data stored in their leaves:

```
data Bin a = Leaf a | Bin (Bin a) (Bin a)
```

A value of type `Bin a` consists of a choice between two constructors. For the first choice, it also contains a value of type `a` whereas for the second it contains two subtrees as children. This means that the `Bin a` type is isomorphic to `Either a (Bin a, Bin a)`. Different libraries differ on how they define their underlying representations. The representation of `Bin a` in terms of *pattern functors* is written as:

$$\text{Rep (Bin } a) = K1\ R\ a\ :+: (K1\ R\ (\text{Bin } a) \ :*:\ K1\ R\ (\text{Bin } a))$$

The `Rep (Bin a)` above is a direct translation of `Either a (Bin a, Bin a)`, but using the combinators provided by `GHC.Generics`. In addition, we also have two conversion functions `from :: a → Rep a` and `to :: Rep a → a` which form an isomorphism between `Bin a` and `Rep (Bin a)`. The interface ties everything under a typeclass:

```
class Generic a where
  type Rep a :: *
  from :: a      → Rep a
  to   :: Rep a → a
```

```

size (Bin (Leaf 1) (Leaf 2))
= gsize (fromgen (Bin (Leaf 1) (Leaf 2)))
= gsize (R1 (K1 (Leaf 1) :+: K1 (Leaf 2)))
= gsize (K1 (Leaf 1)) + gsize (K1 (Leaf 2))
= size (Leaf 1) + size (Leaf 2)
= gsize (fromgen (Leaf 1)) + gsize (fromgen (Leaf 2))
= gsize (L1 (K1 1)) + gsize (L1 (K1 2))
= size (1 :: Int) + size (2 :: Int)

```

FIGURE 2.9: Reduction of $\text{size } (\text{Bin } (\text{Leaf } 1) (\text{Leaf } 2))$.

Defining a generic function is done in two steps. First, we define a class that exposes the function for arbitrary types, in our case, *size*, which we implement for any type via *gsiz*e:

```

class Size (a :: *) where
  size :: a → Int
instance (Size a) ⇒ Size (Bin a) where
  size = gsize ∘ fromgen

```

Next we define the *gsiz*e function that operates on the level of the representation of datatypes. We have to use another class and the instance mechanism to encode a definition by induction on representations:

```

class GSize (rep :: * → *) where
  gsize :: rep x → Int
instance (GSize f, GSize g) ⇒ GSize (f :+: g) where
  gsize (f :+: g) = gsize f + gsize g
instance (GSize f, GSize g) ⇒ GSize (f :+ g) where
  gsize (L1 f) = gsize f
  gsize (R1 g) = gsize g

```

We still have to handle the cases where we might have an arbitrary type in a position, modeled by the constant functor *K1*. We require an instance of *Size* so we can successfully tie the recursive knot.

```

instance (Size x) ⇒ GSize (K1 R x) where
  gsize (K1 x) = size x

```

To finish the description of the generic *size*, we also need instances for the *unit*, *void* and *metadata* pattern functors, called *U1*, *V1*, and *M1* respectively. Their *GSize* is rather uninteresting, so we omit them for the sake of conciseness.

This technique of *mutually recursive classes* is quite specific to the `GHC.Generics` flavor of generic programming. Figure 2.9 illustrates how the compiler goes about choosing instances for computing *size* (*Bin* (*Leaf* 1) (*Leaf* 2)). In the end, we just need an instance for *Size Int* to compute the final result. Literals of type *Int* illustrate what we often call *opaque types*: those types that constitute the base of the universe and are *opaque* to the representation language. This approach to generic programming was also used as the basis for bringing generic programming to Clean [4], where it was later improved to support dynamic types and compile-time optimizations, for example.

2.2.2 EXPLICIT SUMS OF PRODUCTS

The other side of the coin is restricting the shape of the generic values to follow a *sums-of-products* format. This was first done by Löh and de Vries[26] in the `generics-sop` library. The main difference is in the introduction of *Codes*, that limit the structure of representations. If we had access to a representation of the *sum-of-products* structure of *Bin*, we could have defined our *gsize* function following an informal description: sum up the sizes of the fields inside a value, ignoring the constructor.

Unlike `GHC.Generics`, the representation of values is defined by induction on the *code* of a datatype, this code being a type-level list of lists of kind `*`, whose semantics is consonant to a formula in disjunctive normal form. The outer list is interpreted as a sum and each of the inner lists as a product. This section provides an overview of `generics-sop` as required to understand the techniques we use in Chapter 3. We refer the reader to the original paper [26] for a more comprehensive explanation.

Using a *sum-of-products* approach one could write the same *gsize* function shown in Section 2.2.1 as easily as:

$$\begin{aligned} gsize &:: (Generic_{sop} a) \Rightarrow a \rightarrow Int \\ gsize &= sum \circ elim \ (map \ size) \circ from_{sop} \end{aligned}$$

Ignoring the details of *gsize* for a moment, let us focus just on its high level structure. Remembering that *from* now returns a *sum-of-products* view over the data, we are using an eliminator, *elim*, to apply a function to the fields of the constructor used to create a value of type *a*. This eliminator then applies *map size* to the fields of the constructor, returning something akin to a `[Int]`. We then *sum* them up to obtain the final size.

Codes consist of a type-level list of lists. The outer list represents the constructors of a type, and will be interpreted as a sum, whereas the inner lists are interpreted as the fields of the respective constructors, interpreted as products. The `'` sign in the code

below marks the list as operating at the type-level, as opposed to term-level lists which exist at run-time. This is an example of Haskell’s *datatype* promotion [113].

```
type family   Codesop (a :: *) :: '[*]
type instance Codesop (Bin a) = '['[a], '[Bin a, Bin a]]
```

The *representation* is then defined by induction on *Code_{sop}* by the means of generalized *n*-ary sums, *NS*, and *n*-ary products, *NP*. With a slight abuse of notation, one can view *NS* and *NP* through the lens of the following type isomorphisms:

$$\begin{aligned} NS\ f\ [k_1, k_2, \dots] &\equiv f\ k_1\ :\!+\!:\ (f\ k_2\ :\!+\!:\ \dots) \\ NP\ f\ [k_1, k_2, \dots] &\equiv f\ k_1\ :\!*\!:\ (f\ k_2\ :\!*\!:\ \dots) \end{aligned}$$

If we define *Rep_{sop}* to be *NS* (*NP* (*K1 R*)), where **data** *K1 R a* = *K1 a* is borrowed from *GHC.Generics*, we get exactly the representation that *GHC.Generics* issues for *Bin a*. Nevertheless, note how we already need the parameter *f* to pass *NP* to *NS* here.

$$\begin{aligned} Rep_{sop}\ (Bin\ a) &\equiv NS\ (NP\ (K1\ R))\ (Code_{sop}\ (Bin\ a)) \\ &\equiv K1\ R\ a\ :\!+\!:\ (K1\ R\ (Bin\ a)\ :\!*\!:\ K1\ R\ (Bin\ a)) \\ &\equiv Rep_{gen}\ (Bin\ a) \end{aligned}$$

It makes no sense to go through the trouble of adding the explicit *sums-of-products* structure to forget this information in the representation. Instead of piggybacking on *pattern functors*, we define *NS* and *NP* from scratch using *GADTs* [111]. By pattern matching on the values of *NS* and *NP* we inform the type checker of the structure of *Code_{sop}*.

```
data NS :: (k → *) → [k] → * where
  Here  :: f k      → NS f (k ' : ks)
  There :: NS f ks → NS f (k ' : ks)

data NP :: (k → *) → [k] → * where
  ε     ::          NP f []
  (×) :: f x → NP f xs → NP f (x ' : xs)
```

Finally, since our atoms are of kind ***, we can use the identity functor, *I*, to interpret those and define the final representation of values of a type *a* under the *SOP* view:

```
type Repsop a = NS (NP I) (Codesop a)
newtype I (a :: *) = I {unI :: a}
```

To support the claim that one can define general combinators for working with these representations, let us look at *elim* and *map*, used to implement the *gsize* function in the


```

gsize :: (Genericsop a, All2 Size (Codesop a)) ⇒ a → Int
gsize = sum ∘ hcollapse
      ∘ hmap (Proxy :: Proxy Size) (mapIK size) ∘ fromsop

```

FIGURE 2.10: Definition of *gsize* in the *generics-sop* style.

beginning of the section. The *elim* function just drops the constructor index and applies *f*, whereas the *map* applies *f* to all elements of a product.

```

elim :: (∀ k . f k → a) → NS f ks → a
elim f (Here x) = f x
elim f (There x) = elim f x
map :: (∀ k . f k → a) → NP f ks → [a]
map f ε         = []
map f (x × xs) = f x : map f xs

```

Reflecting on the current definition of *size* and comparing it to the `GHC.Generics` implementation of *size*, we see two improvements: (A) we need one fewer typeclass, *GSize*, and, (B) the definition is combinator-based. Considering that the generated *pattern functor* representation of a Haskell datatype will already be in a *sums-of-products*, we do not lose anything by enforcing this structure.

There are still downsides to this approach. A notable one is the need to carry constraints around: the actual *gsize* written with the `generics-sop` library and no sugar is shown in Figure 2.10.

Where *hcollapse* and *hmap* are analogous to the *elim* and *map* combinators defined above. The *All2 Size (Code_{sop} a)* constraint tells the compiler that all of the types serving as atoms for *Code_{sop} a* are an instance of *Size*. Here, *All2 Size (Code_{sop} (Bin a))* expands to *(Size a, Size (Bin a))*. The *Size* constraint also has to be passed around with a *Proxy* for the eliminator of the *n*-ary sum. This is a direct consequence of a *shallow* encoding: since we only unfold one layer of recursion at a time, we have to carry proofs that the recursive arguments can also be translated to a generic representation. We can relieve this burden by recording, explicitly, which fields of a constructor are recursive or not, which is exactly how we start to shape `generics-mrsop` in Chapter 3.

2.2.3 DISCUSSION

Most other generic programming libraries follow a similar pattern of defining the *description* of a datatype in the provided uniform language by some type-level information, and two functions witnessing an isomorphism. The most important feature of such a library

	Pattern Functors	Codes
No Explicit Recursion	<code>GHC.Generics</code>	<code>generics-sop</code>
Simple Recursion	<code>regular</code>	
Mutual Recursion	<code>multirec</code>	

FIGURE 2.11: *Spectrum of static generic programming libraries.*

is how this description is encoded and which primitive operations are used for constructing such encodings. Some libraries, mainly deriving from the SYB approach [53, 77], use the *Data* and *Typeable* typeclasses instead of static type-level information to provide generic functionality – these are a completely different strand of work from what we seek. The main approaches that rely on type-level representations of datatypes are shown in Figure 2.11. These can be compared in their treatment of recursion and on their choice of type-level combinators used to represent generic values.

RECURSION STYLE. There are two ways to define the representation of values. Either we place explicit information about which fields of the constructors of the datatype in question are recursive or we do not.

If we do not mark recursion explicitly, *shallow* encodings are the easier option, where only one layer of the value is turned into a generic form by a call to *from*. This is the kind of representation we get from `GHC.Generics`. The other side of the spectrum would be the *deep* representation, in which the entire value is turned into the representation that the generic library provides in one go.

Marking the recursion explicitly, like in `regular` [82], allows one to choose between *shallow* and *deep* encodings at will. These representations are usually more involved as they need an extra mechanism to represent recursion. In the *Bin* example, the description of the *Bin* constructor changes from “this constructor has two fields of the *Bin a* type” to “this constructor has two fields in which you recurse”. Therefore, a *deep* encoding requires some explicit *least fixpoint* combinator – usually called *Fix* in Haskell.

Depending on the use case, a shallow representation might be more efficient if only part of the value needs to be inspected. On the other hand, deep representations are sometimes easier to use, since the conversion is performed in one go, and afterwards one only has to work with the constructs from the generic library.

The fact that we mark explicitly when recursion takes place in a datatype gives some additional insight into the description. Some functions really need the information about which fields of a constructor are recursive and which are not, like the generic *map* and the generic *Zipper*. This additional power has also been used to define regular expressions over Haskell datatypes [97], for example.

PATTERN FUNCTORS VERSUS CODES. Most generic programming libraries build their type-level descriptions out of three basic combinators: (1) *constants*, which indicate a type is atomic and should not be expanded further; (2) *products* (usually written as `:*`) which are used to build tuples; and (3) *sums* (usually written as `:+:`) which encode the choice between constructors. The *Rep* (*Bin a*) shown before is expressed in this form. Note, however, that there is no restriction on *how* these can be combined. These combinators are usually referred to as *pattern functors*. The *pattern functor*-based libraries are too permissive though, for instance, *K1 R Int* `:*` *Maybe* is a perfectly valid `GHC.Generics pattern functor` but will break generic functions, i.e., *Maybe* is not a combinator.

In practice, one can always use a sum of products to represent a datatype – a sum to express the choice of constructor, and within each constructor a product to declare which fields you have. The `generic-sop` library [26] explicitly uses a list of lists of types, the outer one representing the sum and each inner one thought of as products.

$$Code_{\text{sop}} (Bin\ a) = '['[a], '[Bin\ a, Bin\ a]]$$

The shape of this description follows more closely the shape of Haskell datatypes, and make it easier to implement generic functionality.

Note how the *codes* are different from the *representation*, the latter being defined by induction on the former. This is quite a subtle point and it is common to see both terms being used interchangeably. Here, the *representation* is mapping the *codes*, of kind `'['*]`, into `*`. The *code* can be seen as the format that the *representation* must adhere to. Previously, in the pattern functor approach, the *representation* was not guaranteed to have a certain structure. The expressivity of the language of *codes* is proportional to the expressivity of the combinators the library can provide.



GENERIC PROGRAMMING WITH MUTUALLY RECURSIVE TYPES

The syntax of many programming languages is expressed through a mutually recursive family of datatypes. Before writing a generic differencing algorithm we need to be able to program generically over mutually recursive families of datatypes. Consider Haskell itself, a **do** block constructs an expression, even though the **do** block itself is composed by a list of statements which may include expressions.

```
data Expr = ... | Do [Stmt] | ...  
data Stmt = Assign Var Expr | Let Var Expr
```

Another example is found in HTML and XML documents. These are easily described by a Rose tree, which albeit being a nested type [17], is naturally encoded in the mutually recursive family of datatypes below.

```
data Rose a = Fork a [Rose a]  
data [] a = [] | a : [a]
```

Working with generic mutually recursive families in Haskell, however, is a non-trivial task. The best solution at the time of writing is the `multirec` [112] library, which is unfortunately unfit for writing complex programs – the lack of a combinator-based approach to generic programming and the pattern functor (Section 2.2.1) approach makes it hard to write involved algorithms.

This meant we had to engineer new generic programming libraries to tackle the added complexity of mutual recursion. We have devised two different ways of doing so. First, we wrote the `generics-mrsop` [75] library, which combines a combinator based (Section 2.2.2) approach to generic programming with mutually recursive types. In fact, `generics-mrsop` lies in the intersection of `multirec` and the more modern `generics-sop` [26]. It is worth noting that neither of the aforementioned libraries *compete* with our work. We extend both in orthogonal directions, resulting in a new design altogether, that takes advantage of some modern Haskell extensions which the authors of the previous work could not employ.

The `generics-mrsop` library, Section 3.1, was a conceptual success. It enabled us to prototype and tweak the algorithms discussed in Chapter 4 and Chapter 5 with ease. Yet, a memory leak in the Glasgow Haskell Compiler¹ made it unusable for encoding real programming languages such as those in the `language-python` or `language-java` packages. This frustrating outcome meant that a different approach – which did not rely as heavily on type families – was necessary to look at real-world software version control conflict data.

It turns out that we can sacrifice the sums-of-products structure of `generics-mrsop`, significantly decreasing the reliance of type families, while maintaining a combinator-based approach. This would still enable us to write the algorithms underlying the `hdiff` tool (Chapter 5). This lead us to develop the `generics-simplistic` library, Section 3.2, that still maintains a list of the types that belong in the family, but does not record their internal sum-of-products structure.

This chapter, then, is concerned with explaining our work extending the existing generic programming capabilities of Haskell to support mutually recursive types. We introduce two conceptually different approaches, but with similar expressivity. In Section 3.1 we explore the `generics-mrsop` library. With its ability of representing explicit sums of products we are able to illustrate the `gdif` [55] differencing algorithm, which follows the classical tree-edit distance but in a typed fashion. Then, in Section 3.2, we explore the `generics-simplistic` library, which works on the pattern functor spectrum of generic programming.

3.1 THE GENERICS-MRSOP LIBRARY

The `generics-mrsop` library is an intersection of the `multirec` and `generics-sop` libraries. It uses explicit codes in the sums of products style to guide the representation of datatypes. This enables a simple explicit fixpoint construction and a variety of recursion schemes, which makes the development of generic programs fairly straightforward.

¹<https://gitlab.haskell.org/ghc/ghc/issues/17223> and <https://gitlab.haskell.org/ghc/ghc/issues/14987>

3.1.1 EXPLICIT FIXPOINTS WITH CODES

Introducing information about the recursive positions in a type requires more expressive codes than in Section 2.2.2. Where our *codes* were a list of lists of types, which could be anything, we now have a list of lists of *Atom*, which maintains information about whether a position is recursive or not.

```

data Atom = I | KInt | ...
type family Codefix (a :: *) :: '[Atom]
type instance Codefix (Bin Int) = '[KInt], '[I, I]

```

Here, *I* is used to mark the recursive positions and *KInt*, ... are codes for a predetermined selection of primitive types, which we refer to as *opaque types*. Favoring the simplicity of the presentation, we will stick with only hard coded *Int* as the only opaque type in the universe. Later on, in Section 3.1.2.1, we parameterize the whole development by the choice of opaque types.

We can no longer represent polymorphic types in this universe – the *codes* themselves are not polymorphic. Back in Section 2.2.2 we have defined *Code_{sop}* (*Bin a*), and this would work for any *a*. The lack of polymorphism might seem like a disadvantage at first, but if we are interested in deep generic representations, it is actually an advantage, as it allows us to have a deep conversion for free as we do not need to carry *Generic* constraints around. That is, say we want to deeply convert a value of type *Bin a* to its generic representation polymorphically on *a*. We can only do so if we have access to the *Code_{sop}* *a*, which comes from knowing *Generic a*. By specifying the types involved beforehand, we are able to get by without having to carry all of the constraints we needed in, for instance, *gsiz*e at the end of Section 2.2.2. The main benefit is in the simplicity of combinators we will define in Section 3.1.2.2.

The *Rep_{fix}* datatype is similar to the *Rep_{sop}*, but uses an additional layer that maps an *Atom* into ***, denoted *NA*. Since an atom can be either an opaque type, known statically, or some type that must be placed in a recursive position later on, we need just one parameter in *NA*.

```

data NA :: * → Atom → * where
  NAI :: x → NA x I
  NAK :: Int → NA x KInt
newtype Repfix a x = Rep {unRep :: NS (NP (NA x)) (Codefix a)}

```

The *Generic_{fix}* typeclass, below, witnesses the isomorphism between ordinary types and their deep sums-of-products representation. Similarly to the other generic type-classes out there, it contains just the familiar *to_{fix}* and *from_{fix}* components. We illustrate

part of the instance that witnesses that *Bin Int* has a generic representation below. We omit the to_{fix} function as it is the opposite of $from_{fix}$.

```

class Genericfix a where
  fromfix :: a → Repfix a a
  tofix   :: Repfix a a → a

instance Genericfix (Bin Int) where
  fromfix (Leaf x) = Rep (      Here (NAK x × ε))
  fromfix (Bin l r) = Rep (There (Here (NAI l × NAI r × ε)))

```

It is an interesting exercise to implement the *Functor* instance for $(Rep_{fix} a)$ – where it can be seen that we were only able to lift it to a functor by recording the information about the recursive positions. Otherwise, there would be no easy way of knowing where to apply f when defining $fmap f$.

Nevertheless, working directly with Rep_{fix} is hard – we need to pattern match on *There* and *Here*, whereas we actually want to have the notion of *constructor* for the generic setting too! The main advantage of the *sum-of-products* structure is to allow a user to pattern match on generic representations just like they would on values of the original type, contrasting with `GHC.Generics`. One can precisely state that a value of a representation is composed by a choice of constructor and its respective product of fields by the *View* type. This *view* pattern [108, 66] is common in dependently typed programming.

```

data Nat = Z | S Nat
data View :: [[Atom]] → * → * where
  Tag :: Constr n t → NP (NA x) (Lkup t n) → View t x

```

A value of *Constr n sum* is a proof that n is a valid constructor for *sum*, stating that $n < length\ sum$. *Lkup* performs list lookup at the type-level. To improve type error messages, we generate a *TypeError* whenever we reach a given index n that is out of bounds. Interestingly, our design guarantees that this case is never reached by *Constr*.

```

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

type family Lkup (ls :: [k]) (n :: Nat) :: k where
  Lkup '[]      _      = TypeError "Index out of bounds"
  Lkup (x : xs) 'Z      = x
  Lkup (x : xs) ('S n) = Lkup xs n

```

With the help of *sop* and *inj*, declared below, we are able to pattern match and inject into generic values. Unfortunately, matching on *Tag* directly can be cumbersome, but

```

crush :: (Genericfix a) ⇒ (∀ x . Int → b) → ([b] → b) → a → b
crush k cat = crushFix ∘ deepFrom
  where crushFix :: Fix (Repfix a) → b
        crushFix = cat ∘ elimNS (elimNP go) ∘ unFix
        go (NAI x) = crushFix x
        go (NAK i) = k i
    
```

 FIGURE 3.1: *Generic crush combinator.*

we can always use pattern synonyms [90] to circumvent that. For example, the synonyms below describe the constructors *Bin* and *Leaf*.

```

pattern (Pat Leaf) x = Tag CZ (NAK x × ε)
pattern (Pat Bin) l r = Tag (CS CZ) (NAI l × NAI r × ε)
inj :: View sop x → Repfix sop x
sop :: Repfix sop x → View sop x
    
```

Having the core of the *sums-of-products* universe defined, we turn our attention to the representation of recursion through the *Fix* datatype. This enables us to convert values to their *deep* representation.

CONVERTING TO A DEEP REPRESENTATION. The $from_{fix}$ function still returns a shallow representation. But by constructing the least fixpoint of $Rep_{fix} a$ we can easily obtain the deep encoding for free, by recursively translating each layer of the shallow encoding.

```

newtype Fix f = Fix {unFix :: f (Fix f)}
deepFrom :: (Genericfix a) ⇒ a → Fix (Repfix a)
deepFrom = Fix ∘ fmap deepFrom ∘ fromfix
    
```

So far, we handle the same class of types as the `regular` [82] library, but we require the representation to follow a sums of products structure by the means of $Code_{fix}$. Those types are guaranteed to have an initial algebra, and indeed, the generic catamorphism is defined as expected:

```

fold :: (Repfix a b → b) → Fix (Repfix a) → b
fold f = f ∘ fmap (fold f) ∘ unFix
    
```

Some functions may consume a value and produce a single value, but do not need the full expressivity of *fold*. Instead, if we know how to consume the opaque types and combine those results, we can consume any $Generic_{fix}$ type using *crush*, which is defined

in Figure 3.1. The behavior of *crush* is defined by (1) how to turn atoms into the output type *b* – in this case we only have integer atoms, and thus we require an *Int* \rightarrow *b* function – and (2) how to combine the values bubbling up from each member of a product.

Finally, we come full circle to our running *gsize* example as it was promised in the introduction. This is noticeably the smallest implementation so far, and very straight to the point.

```
gsize :: (Genericfix a)  $\Rightarrow$  a  $\rightarrow$  Int
gsize = crush (const 1) sum
```

At this point we have combined the insight from the *regular* library of keeping track of recursive positions with the convenience of the *generics-sop* for enforcing a specific *normal form* on representations. By doing so, we were able to provide a deep encoding for free. This essentially frees us from the burden of maintaining constraints. Compare the *gsize* above with its *generics-sop* variation, in Figure 2.10. The information about the recursive position allows us to write concise combinators, such as *crush*, and a convenient *View* type for easy generic pattern matching. The only thing keeping us from handling larger applications is the limited form of recursion.

3.1.2 MUTUAL RECURSION

Conceptually, going from regular types (Section 3.1.1) to mutually recursive families is simple. We just need to reference not only one type variable, but one for each element in the family. This is usually [58, 5] done by adding an index to the recursive positions to represents each member of the family. As a running example, we use the familiar *rose tree* family.

```
data Rose a = Fork a [Rose a]
data [] a = [] | a : [a]
```

The previously introduced *Code_{fix}*, Section 3.1.1, is not expressive enough to describe this datatype. In particular, when we try to write *Code_{fix}* (*Rose Int*), there is no immediately recursive appearance of *Rose* itself, so we cannot use the atom *I* in that position. Furthermore [*Rose a*] is not an opaque type either, so we cannot use any of the other combinators provided by *Atom*. We would like to record information about *Rose Int* referring to itself via another datatype.

Our solution is to move from codes of datatypes to *codes for families of datatypes*. We no longer talk about *Code_{fix}* (*Rose Int*) or *Code_{fix}* [*Rose Int*] in isolation. Codes only make sense within a family, that is, a list of types. Hence, we talk about the codes of the two types in the family: *Code_{mrec}* '[*Rose Int*, [*Rose Int*]]'. Then we extend the language

of *Atoms* by appending to *I* a natural number which specifies the member of the family to recurse into:

```
data Atom = I Nat | KInt | ...
```

The code of this recursive family of datatypes can be described as:

```
type FamRose      = '[Rose Int, [Rose Int]]
type Codemrec FamRose = '[ '[KInt, I (S Z)]
                          , '[[], '[I Z, I (S Z)]]
                          ]
```

Let us have a closer look at the code for *Rose Int*, which appears in the first place in the list. There is only one constructor which has an *Int* field, represented by *KInt*, and another in which we recurse via the second member of our family (since lists are 0-indexed, we represent this by *S Z*). Similarly, the second constructor of *[Rose Int]* points back to both *Rose Int* using *I Z* and to *[Rose Int]* itself via *I (S Z)*.

Having settled on the definition of *Atom*, we now need to adapt *NA* to the new *Atoms*. To interpret any *Atom* into $*$, we need a way assign values to the different recursive positions. This information is given by an additional type parameter φ that maps natural numbers into types.

```
data NA :: (Nat → *) → Atom → * where
  NAI  ::  $\varphi$  n → NA  $\varphi$  (I n)
  NAK  :: Int  → NA  $\varphi$  KInt
```

This additional φ naturally bubbles up to *Rep_{mrec}*.

```
type Repmrec ( $\varphi$  :: Nat → *) (c :: [[Atom]]) = NS (NP (NA  $\varphi$ )) c
```

The only piece missing here is tying the recursive knot. If we want our representation to describe a family of datatypes, the obvious choice for φ *n* is to look up the type at index *n* in *FamRose*. In fact, we are simply performing a type-level lookup in the family, so we can reuse the *Lkup* from Section 3.1.1.

In principle, this is enough to provide a ground representation for the family of types. Let *fam* be a family of types, like *'[Rose Int, [Rose Int]]*, and *codes* the corresponding list of codes. Then the representation of the type at index *ix* in the list *fam* is given by:

```
Repmrec (Lkup fam) (Lkup codes ix)
```

This definition states that to obtain the representation of the type at index ix , we first lookup its code. Then, in the recursive positions we interpret each $I\ n$ by looking up the type at that index in the original family. This gives us a *shallow* representation.

Unfortunately, Haskell only allows saturated, that is, fully-applied type families. Hence, we cannot partially apply *Lkup* like we did it in the example above. As a result, we need to introduce an intermediate datatype *El*,

```
data El :: [*] → Nat → * where
  El :: Lkup fam ix → El fam ix
```

The representation of the family *fam* at index ix is thus given in terms of *El*, which can be partially applied, $Rep_{mrec} (El\ fam) (Lkup\ codes\ ix)$. We only need to use *El* in the first argument, because that is the position in which we require partial application. The second position has *Lkup* already fully-applied, and can stay as is.

We still have to relate a family of types to their respective codes. As in other generic programming approaches, we want to make their relation explicit. The *Family* typeclass below realizes this relation, and introduces functions to perform the conversion between our representation and the actual types. Using *El* here spares us from using a proxy for *fam* in $from_{mrec}$ and to_{mrec} :

```
class Family (fam :: [*]) (codes :: [[[Atom]]]) where
   $from_{mrec} :: SNat\ ix \rightarrow El\ fam\ ix \rightarrow Rep_{mrec}\ (El\ fam)\ (Lkup\ codes\ ix)$ 
   $to_{mrec} \quad :: SNat\ ix \rightarrow Rep_{mrec}\ (El\ fam)\ (Lkup\ codes\ ix) \rightarrow El\ fam\ ix$ 
```

One of the differences between other approaches and ours is that we do not use an associated type to define the *codes* for the family *fam*. This path is that it alleviates the burden of writing the longer $Code_{mrec}\ fam$ every time we want to refer to *codes*. Furthermore, there are types like lists which appear in many different families, and in that case it makes sense to speak about a relation instead of a function.

Since $from_{mrec}$ and to_{mrec} operate on families, we have to specify how to translate *each* of the members of the family *to* and *from* their generic representation. This translation needs to know the index of the datatype we are converting between in each case, hence the additional singleton *SNat* ix parameter. Pattern matching on this singleton [29] type informs the compiler about the shape of the *Nat* index. Its definition is:

```
data SNat (n :: Nat) where
  SZ :: SNat 'Z
  SS :: SNat n → SNat ('S n)
```

The *SNat* datatype, in turn, enables us to write the definition of $from_{mrec}$ for the mutually recursive family of rose trees.

```
-- First type in the family
from_mrec SZ (El (Fork x ch)) = Rep (Here (NA_K x × NA_I ch × ε))
-- Second type in the family
from_mrec (SS SZ) (El [])      = Rep (Here ε)
from_mrec (SS SZ) (El (x : xs)) = Rep (There (Here (NA_I x × NA_I xs × ε)))
```

By pattern matching on the index, the compiler knows which family member to expect as a second argument. This then allows the pattern matching on the *El* to typecheck.

The limitations of the Haskell type system lead us to introduce *El* as an intermediate datatype. Our $from_{mrec}$ function does not take a member of the family directly, but an *El*-wrapped one. However, to construct that value, *El* needs to know its parameters, which amounts to knowing the family we are embedding our type into and the index in that family. Those values are not immediately obvious, but we can use Haskell’s visible type application [30] to work around it. The *into* function injects a value into the corresponding *El*:

```
into :: ∀ fam ty ix . (ix ~ Idx ty fam , Lkup fam ix ~ ty) ⇒ ty → El fam ix
into = El

intoRose :: Rose Int → El FamRose 'Z
intoRose = into @FamRose
```

Idx, here, is a closed type family implementing the inverse of *Lkup*, that is, obtaining the index of the type *ty* in the list *fam*. Using this function we can turn a $[Rose\ Int]$ into its generic representation by writing $from_{mrec} \circ into\ @FamRose$. The type application $@FamRose$ is responsible for fixing the mutually recursive family we are working with, which allows the type checker to reduce all the constraints and happily inject the element into *El*.

DEEP REPRESENTATION. In Section 3.1.1 we have described a technique to derive deep representations from shallow representations. We can play a very similar trick here. The main difference is the definition of the least fixpoint combinator, which receives an extra parameter of kind *Nat* indicating which *code* to use first:

```
newtype Fix (codes :: [[Atom]]) (ix :: Nat)
  = Fix {unFix :: Rep_mrec (Fix codes) (Lkup codes ix)}
```

Intuitively, since now we can recurse on different positions, we need to keep track of the representations for all those positions in the type. This is the job of the *codes* argument. Furthermore, our *Fix* does not represent a single datatype, but rather the

whole family. Thus, we need each value to have an additional index to declare on which element of the family it operates.

As in the previous section, we can obtain the deep representation by iteratively applying the shallow representation. Earlier we used *fmap* since the *Rep_{fix}* type was a functor. *Rep_{mrec}* on the other hand cannot be given a *Functor* instance, but we can still define a similar function *mapRec*,

$$\text{mapRep} :: (\forall ix . \varphi_1 ix \rightarrow \varphi_2 ix) \rightarrow \text{Rep}_{\text{mrec}} \varphi_1 c \rightarrow \text{Rep}_{\text{mrec}} \varphi_2 c$$

This signature tells us that if we want to change the φ_1 argument in the representation, we need to provide a natural transformation from φ_1 to φ_2 , that is, a function which works over each possible index this φ_1 can take and does not change this index. This follows from φ_1 having kind *Nat* \rightarrow $*$.

$$\begin{aligned} \text{deepFrom} &:: \text{Family fam codes} \Rightarrow \text{El fam ix} \rightarrow \text{Fix} (\text{Rep}_{\text{mrec}} \text{ codes ix}) \\ \text{deepFrom} &= \text{Fix} \circ \text{mapRec} \text{ deepFrom} \circ \text{from}_{\text{mrec}} \end{aligned}$$

ONLY WELL-FORMED REPRESENTATIONS ARE ACCEPTED. At first glance, it may seem like the *Atom* datatype gives too much freedom: its *I* constructor receives a natural number, but there is no apparent static check that this number refers to an actual member of the recursive family we are describing. For example, the list of codes given by `'[[[KInt, I (S (S Z))]]]` is accepted by the compiler although it does not represent any family of datatypes.

A direct solution to this problem is to introduce yet another index, this time in the *Atom* datatype, which specifies which indices are allowed. The *I* constructor is then refined to take not any natural number, but only those which lie in the range – this is usually known as *Fin n*.

```
data Atom (n :: Nat) = I (Fin n) | KInt | ...
```

The lack of dependent types makes this approach very hard, in Haskell. We would need to carry around the inhabitants *Fin n* and define functionality to manipulate them, which would greatly hinder the usability of the library.

By looking a bit more closely, we find that we are not losing any type-safety by allowing codes which reference an arbitrary number of recursive positions. Users of our library are allowed to write the previous ill-defined code, but when trying to write *values* of the representation of that code, the *Lkup* function detects the out-of-bounds index, raising a type error and preventing the program from compiling in the first place, instead of crashing at run-time.

3.1.2.1 PARAMETERIZED OPAQUE TYPES

Up to this point we have considered *Atom* to include a predetermined selection of *opaque types*, such as *Int*, each of them represented by one of the constructors other than *I*. This is far from ideal, for two conflicting reasons:

- a) The choice of opaque types might be too narrow. For example, the user of our library may decide to use *ByteString* in their datatypes. Since that type is not covered by *Atom*, nor by our generic approach, this implies that *generics-mrsop* becomes useless to them.
- b) The choice of opaque types might be too wide. If we try to encompass any possible situation, we end up with a huge *Atom* type. But for a specific use case, we might be interested only in *Ints* and *Floats*. This means we have to worry about possibly ill-formed representations and pattern matches which should never be reached.

Our solution is to *parameterize Atom*, giving users the choice of opaque types. For example, if we only want to deal with numeric opaque types, we can write:

```
data Atom kon    = I Nat | K kon
data NumericK    = KInt | KInteger | KFloat
type NumericAtom = Atom NumericK
```

The representation of codes must be updated to reflect the possibility of choosing different sets of opaque types. The *NA* datatype in this final implementation provides two constructors, one per constructor in *Atom*. The *NS* and *NP* datatypes do not require any change.

```
data NA :: (kon → *) → (Nat → *) → Atom kon → * where
  NAI :: φ n → NA κ φ (I n)
  NAK :: κ k → NA κ φ (K k)
type Repmrec (κ :: kon → *) (φ :: Nat → *) (c :: [[Atom kon]]) = NS (NP (NA κ φ)) c
```

The *NA_K* constructor in *NA* makes use of an additional argument κ . The problem is that we are defining the code for the set of opaque types by a specific kind, such as *Numeric* above. On the other hand, values which appear in a field must have a type whose kind is $*$. Thus, we require a mapping from each of the codes to the actual opaque type they represent. This is exactly the *opaque type interpretation* κ . Here is the datatype interpreting *NumericK* into ground types:

```
data NumericI :: NumericK → * where
  IInt :: Int → NumericI KInt
  IFloat :: Float → NumericI KFloat
```

```

class Family ( $\kappa :: \text{kon} \rightarrow *$ ) (fam ::  $[*]$ ) (codes ::  $[[[Atom \text{kon}]]]$ ) where
  frommrec ::  $SNat \text{ix} \rightarrow El \text{fam} \text{ix} \rightarrow Rep_{mrec} \kappa (El \text{fam}) (Lkup \text{codes} \text{ix})$ 
  tomrec   ::  $SNat \text{ix} \rightarrow Rep_{mrec} \kappa (El \text{fam}) (Lkup \text{codes} \text{ix}) \rightarrow El \text{fam} \text{ix}$ 

```

FIGURE 3.2: *Family* typeclass with support for different opaque types.

The last piece of our framework which has to be updated to support different sets of opaque types is the *Family* typeclass, as given in Figure 3.2. This typeclass provides an interesting use case for the new dependent features in Haskell; both κ and *codes* are parameterized by an implicit argument *kon* which represents the set of opaque types.

We stress that the parametrization over opaque types does *not* mean that we can use only closed universes of opaque types. It is possible to provide an *open* representation by choosing $(*)$ – the whole kind of Haskell’s ground types – as argument to *Atom*. As a consequence, the interpretation ought to be of kind $* \rightarrow *$, as given by *Value*, below. To use $(*)$ as an argument to a type, we must enable the `TypeInType` language extension [109, 110].

```

data Value ::  $* \rightarrow *$  where
  Value ::  $t \rightarrow Value \ t$ 

```

3.1.2.2 SELECTION OF USEFUL COMBINATORS

The advantages of a *code based* approach to generic programming becomes evident when we look at the generic combinators that `generics-mrsop` provides. We refer the reader to the actual documentation for a comprehensive list. Here we look at a selection of useful functions in their full form. Let us start with the bifunctionality of *Rep_{mrec}*:

```

bimapRep :: ( $\forall k . \kappa_1 k \rightarrow \kappa_2 k$ )  $\rightarrow (\forall ix . \varphi_1 ix \rightarrow \varphi_2 ix)$ 
           $\rightarrow Rep_{mrec} \kappa_1 \varphi_1 c \rightarrow Rep_{mrec} \kappa_2 \varphi_2 c$ 
bimapRep  $f_k f_l = mapNS (mapNP (mapNA f_l f_l))$ 

```

To destruct a *Rep_{mrec}* $\kappa \varphi c$ we need a way for eliminating every recursive position or opaque type inside the representation and a way of combining these results.

```

elimRep :: ( $\forall k . \kappa k \rightarrow a$ )  $\rightarrow (\forall ix . \varphi ix \rightarrow a) \rightarrow ([a] \rightarrow b) \rightarrow Rep_{mrec} \kappa \varphi c \rightarrow b$ 
elimRep  $f_k f_l \text{ cat} = elimNS \text{ cat} (elimNP (elimNA f_k f_l))$ 

```

Another useful operator, particularly when combined with *bimapRep* is the *zipRep*, that works just like a regular *zip*. Our *zipRep* attempts to put two values of a representa-

```

geq :: (EqHO  $\kappa$ , Family  $\kappa$  fam codes)  $\Rightarrow$  ( $\forall k$  .  $\kappa k \rightarrow \kappa k \rightarrow \text{Bool}$ )
     $\rightarrow$  El fam ix  $\rightarrow$  El fam ix  $\rightarrow \text{Bool}$ 
geq eqK x y = go (deepFrom x) (deepFrom y)
  where go (Fix x) (Fix y)
        = maybe False (elimRep (uncurry eqK) (uncurry go) and) $ zipRep x y

```

FIGURE 3.3: *Generic equality.*

tion “side-by-side”, as long as they are constructed with the same injection into the n -ary sum, *NS*.

```

zipRep :: Repmrec  $\kappa_1$   $\varphi_1$  c  $\rightarrow$  Repmrec  $\kappa_2$   $\varphi_2$  c  $\rightarrow$  Maybe (Repmrec ( $\kappa_1$   $∶∶$   $\kappa_2$ ) ( $\varphi_1$   $∶∶$   $\varphi_2$ ) c)
zipRep r s = case (sop r, sop s) of
  (Tag cr pr, Tag cs ps)  $\rightarrow$  case testEquality cr pr of
    Just Refl  $\rightarrow$  inj cr  $\langle\!\langle$  zipWithNP zipAtom pr ps

```

We use *testEquality* from *Data.Type.Equality* to check for type index equality and inform the compiler of that fact by matching on *Refl*.

Finally, we can start assembling these building blocks into more practical functionality. Figure 3.3 shows the definition of generic equality using *generics-mrsop*, where the *EqHO* typeclass is a lifted version of *Eq*, for types of kind $k \rightarrow *$, defined below. The library also provide *ShowHO*, the *Show* counterpart.

```

class EqHO (f :: a  $\rightarrow$  *) where
  eqHO ::  $\forall x$  . f x  $\rightarrow$  f x  $\rightarrow \text{Bool}$ 

```

We decided to provide a custom equality in *generics-mrsop* for two main reasons. Firstly, when we started developing the library the *-XQuantifiedConstraints* [18] extension was not completed. Yet, once quantified constraints were available in Haskell we wrote *generics-mrsop-2.2.0* using the extension and defining *EqHO f* as a synonym to $\forall x \circ \text{Eq} (f x)$. Developing applications on top of *generics-mrsop* became more difficult. The user now would have to reason about and pass around complicated constraints down to datatypes and auxiliary functions. Moreover, our use case was very simple, not extracting any of the advantages of quantified constraints. Eventually we decided to rollback to the lifted *EqHO* presented above in *generics-mrsop-2.3.0*.

As presented so far, we have all the necessary tools to encode our first differencing attempt, shown in Chapter 4 of this thesis. The next sections discusses some aspects that, albeit not directly required for understanding the remainder of this thesis, are interesting in their own right and round off the presentation of *generics-mrsop* as a library.

3.1.3 PRACTICAL FEATURES

The development of the `generics-mrsop` library started primarily to enable us to write `hdiff` (Chapter 5) possible. This was a great expressivity test for our generic programming library and led us to develop overall useful features that, although not novel, make the adoption of a generic programming library much more likely. This section is a small tutorial into two important practical features of `generics-mrsop` and documents the engineering effort that was put in the library.

3.1.3.1 TEMPLATE HASKELL

Having a convenient and robust way to get the *Family* instance for a given selection of datatypes is of paramount importance for the usability of our library. In a real scenario, a mutually recursive family may consist of many datatypes with dozens of constructors. Sometimes these datatypes are written with parameters, or come from external libraries.

Our goal here is to automate the generation of *Family* instances under all those circumstances using *Template Haskell* [100]. From the programmers' point of view, they only need to call *deriveFamily* with the topmost (that is, the first) type of the family. For example:

```
data Exp var = ...
data Stmt var = ...
data Prog var = ...
deriveFamily [t|Prog String|]
```

The *deriveFamily* takes care of unfolding the (type-level) recursion until it reaches a fixpoint. In this case, the type synonym `FamProgString = '[Prog String , ...]` will be generated, together with its *Family* instance. Optionally, one can also pass along a custom function to decide whether a type should be considered opaque. By default, it uses a selection of Haskell built-in types as opaque types.

UNFOLDING THE FAMILY The process of deriving a whole mutually recursive family from a single member is conceptually divided into two disjoint processes. First we repeatedly unfold all definitions and follow all the recursive paths until we reach a fixpoint. At that moment we know that we have discovered all the types in the family. Second, we translate the definition of those types to the format our library expects. During the unfolding process we keep a key-value map in a *State* monad, keeping track of three things: the types we have seen; the types we have seen *and* processed; and the indices of those within the family.

Let us illustrate this process in a bit more detail using our running example of a mutually recursive family and consider what happens within *Template Haskell* when it starts unfolding the *deriveFamily* clause.

```
data Rose a = Fork a [Rose a]
data [a]    = [] | a : [a]
deriveFamily [t|Rose Int|]
```

The first thing that happens is registering that we have seen the type *Rose Int*. Since it is the first type to be discovered, it is assigned index zero within the family. Next we need to reify the definition of *Rose*. At this point, we query *Template Haskell* for the definition, and we obtain **data** *Rose* *x* = *Fork* *x* [*Rose* *x*]. Since *Rose* has kind $* \rightarrow *$, it cannot be directly translated – our library only supports ground types, which are those with kind $*$. But we do not need a generic definition for *Rose*, we just need the specific case where *x* = *Int*. Essentially, we just apply the reified definition of *Rose* to *Int* and β -reduce it, giving us *Fork Int* [*Rose Int*].

The next processing step is looking into the types of the fields of the (single) constructor *Fork*. First we see *Int* and decide it is an opaque type, say *KInt*. Second, we see [*Rose Int*] and notice it is the first time we see this type. Hence, we register it with a fresh index, *S Z* in this case. The final result for *Rose Int* is `'['[K KInt, I (S Z)]]`.

We now go into [*Rose Int*] for processing. Once again we need to perform some amount of β -reduction at the type-level before inspecting its fields. The rest of the process is the same as that for *Rose Int*. However, when we encounter the field of type *Rose Int* this is already registered, so we just need to use the index *Z* in that position.

The final step is generating the actual Haskell code from the data obtained in the previous process. This is a very verbose and mechanical process, whose details we omit. In short, we generate the necessary type synonyms, pattern synonyms, the *Family* instance, and metadata information. The generated type synonyms are named after the topmost type of the family, passed to *deriveFamily*:

```
type FamRoseInt  = '['[Rose Int      , [Rose Int]]
type CodesRoseInt = '['['[K KInt, I (S Z)]] , '['[], '['[I Z, I (S Z)]]]
```

The actual *Family* instance is exactly as the one shown in Section 3.1.2

```
instance Family Singl FamRoseInt CodesRoseInt where ...
```

```

data DatatypeInfo :: [[*]] → * where
  ADT :: ModuleName → DatatypeName → NP ConstrInfo cs → DatatypeInfo cs
  New  :: ModuleName → DatatypeName → ConstrInfo '[c] → DatatypeInfo '['[c]]

data ConstrInfo :: [*] → * where
  Constructor :: ConstrName → ConstrInfo xs
  Infix       :: ConstrName → Associativity → Fixity → ConstrInfo '[x,y]
  Record      :: ConstrName → NP FieldInfo xs → ConstrInfo xs

data FieldInfo :: * → * where
  FieldInfo :: FieldName → FieldInfo a

class HasDatatypeInfo a where
  datatypeInfo :: proxy a → DatatypeInfo (Code a)

```

FIGURE 3.4: Definitions related to metadata from *generics-sop*.

3.1.3.2 METADATA

There is one final ingredient missing to make *generics-mrsop* fully usable in practice. We must maintain the *metadata* information of our datatypes. This metadata includes the datatype name, the module where it was defined, and the name of the constructors. Without this information we would never be able to pretty print the generic code in a satisfactory way. This includes conversion to semi-structured formats, such as JSON, or actual pretty printing.

Like in *generics-sop* [26], having the code for a family of datatypes available allows for a completely separate treatment of metadata. This is yet another advantage of the sum-of-products approach compared to the more traditional pattern functors. In fact, our handling of metadata is heavily inspired from *generics-sop*, so much so that we will start by explaining a simplified version of their handling of metadata, and then outline the differences to our approach.

The general idea is to store the meta information following the sum-of-products structure of the datatype itself. Instead of data, we keep track of the names of the different parts and other meta information that can be useful. It is advantageous to keep metadata separate from the generic representation as it would only clutter the definition of generic functionality. This information is tied to a datatype by means of an additional typeclass *HasDatatypeInfo*. Generic functions may now query the metadata by means of functions like *datatypeName*, which reflect the type information into the term level. The definitions are given in Figure 3.4 and follow closely how *generics-sop* handles metadata.

Our library uses the same approach to handle metadata. In fact, the code remains almost unchanged, except for adapting it to the larger universe of datatypes we can now

handle. Unlike `generic-sop`, our list of lists representing the sum-of-products structure does not contain types of kind `*`, but `Atoms`. All the types representing metadata at the type-level must be updated to reflect this new scenario:

```
data DatatypeInfo :: [[Atom kon]] → * where ...
data ConstrInfo  :: [Atom kon] → * where ...
data FieldInfo   :: Atom kon  → * where ...
```

As we have discussed above, our library is able to generate codes not only for single types of kind `*`, like `Int` or `Bool`, but also for types which are the result of type-level applications, such as `Rose Int` and `[Rose Int]`. The shape of the metadata information in `DatatypeInfo`, a module name plus a datatype name, is not enough to handle these cases. We replace the uses of `ModuleName` and `DatatypeName` in `DatatypeInfo` by a richer promoted type `TypeName`, which can describe applications, as required.

```
data TypeName = ConT ModuleName DatatypeName | TypeName :@: TypeName
data DatatypeInfo :: [[Atom kon]] → * where
  ADT :: TypeName → NP ConstrInfo cs → DatatypeInfo cs
  New :: TypeName → ConstrInfo '[c] → DatatypeInfo '['[c]]
```

An important difference to `generics-sop` is that the metadata is not defined for a single type, but for a type *within* a family. This can be seen in the signature of `datatypeInfo`, which receives proxies for both the family and the type. The type equalities in that signature reflect the fact that the given type `ty` is included with index `ix` within the family `fam`. This step is needed to look up the code for the type in the right position of `codes`.

```
class (Family κ fam codes) ⇒ HasDatatypeInfo κ fam codes ix | fam → κ codes where
  datatypeInfo :: (ix ~ Idx ty fam, Lkup ix fam ~ ty) ⇒ Proxy fam → Proxy ty
               → DatatypeInfo (Lkup ix codes)
```

Template Haskell would generate the instance below for `Rose Int`:

```
instance HasDatatypeInfo Singl FamRose CodesRose Z where
  datatypeInfo _ _ = ADT (ConT "E" "Rose" :@: ConT "Prelude" "Int")
    $ (Constructor "Fork") × ε
```

3.1.4 EXAMPLE: WELL-TYPED CLASSICAL TREE DIFFERENCING

This section, based on the work of Lempink [55] which originally implemented in the `gdif` library, is the related work that is closest to ours in the sense that it is the only

typed approach to differencing. The presentation provided here is adapted from Van Putten’s [92] master thesis and is available as the `generics-mrsop-gdiff` library.

Next, we discuss how to make tree edit-scripts (Section 2.1.2) type-safe following the work of Lempink [55]. We start by lifting edit-scripts to kind $[*] \rightarrow [*] \rightarrow *$, which enables the indexing of the types for the source and destination forests of particular edit-scripts. Consequently, instead of differencing a list of trees, we will difference an n -ary product, NP , indexed by the type of each tree.

```

type PatchGD κ codes xs ys = ES κ codes xs ys
diff :: (TestEquality κ, EqHO κ)
      ⇒ NP (NA κ (Fix κ codes)) xs → NP (NA κ (Fix κ codes)) ys
      → PatchGD κ codes xs ys

```

One confusing complication is that our edit operations operate over both constructors of the family and opaque values, unlike the untyped version of tree differencing (Section 2.1.2), where everything is a label. Consequently, writing the edit operations requires a uniform treatment of recursive constructors and opaque values, which is done by the `Cof` type, read as *constructor-of*. This represents the *unit of modification* of each edit operation. A value of type `Cof κ codes at tys` represents a constructor of atom *at*, which expects arguments whose type is `NP I tys`, for the family *codes* with opaque types interpreted by κ . Its definition is given below.

```

data Cof κ codes :: Atom kon → [Atom kon] → * where
  ConstrI :: (IsNat c, IsNat n)
           ⇒ Constr (Lkup n codes) c → ListPrf (Lkup c (Lkup n codes))
           → Cof κ codes ('I n) (Lkup c (Lkup n codes))
  ConstrK :: κ k → Cof κ codes ('K k) Pnil

```

We need the `ListPrf` argument to `ConstrI` to be able to manipulate the type-level lists when defining the application function, `applyES`. But first, we have to define our edit-scripts. A value of type `ES κ codes xs ys` represents a transformation of a value of `NP (NA κ (Fix κ codes)) xs` into a value of `NP (NA κ (Fix ki codes)) ys`. The `NP` serves as a list of trees, as is usual for the tree differencing algorithms, but it enables us to keep track of the type of each individual tree through the index to `NP`.

```

data ES κ codes :: [Atom kon] → [Atom kon] → * where
  ESO :: ES κ codes '[] '[]
  Ins :: Cof κ codes a t → ES κ codes i (t : # : j) → ES κ codes i (a ' : j)
  Del :: Cof κ codes a t → ES κ codes (t : # : i) j → ES κ codes (a ' : i) j
  Cpy :: Cof κ codes a t → ES κ codes (t : # : i) (t : # : j) → ES κ codes (a ' : i) (a ' : j)

```

Let us take *Ins*, for example. Inserting a constructor $c :: t_1 \rightarrow \dots \rightarrow t_n \rightarrow 'I \text{ ix}$ in a forest $x_1 \times x_2 \times \dots \times \text{Nil}$ will take the first n elements of that forest and use them as arguments to c . This is realized by the *insCof* function, shown below.

```

insCof :: Cof  $\kappa$  codes  $a \ t$ 
         $\rightarrow NP (NA \ \kappa \ (Fix \ \kappa \ codes)) (t : \# : xs) \rightarrow NP (NA \ \kappa \ (Fix \ \kappa \ codes)) (a' : xs)$ 
insCof (ConstrK  $k$ )       $xs = NA_K \ k \times xs$ 
insCof (ConstrI  $c \ ispoa$ )  $xs = \text{let } (poa, xs') = \text{split } ispoa \ xs \text{ in } NA_I (Fix \$ \text{inj } c \ poa) \times xs'$ 

```

The example also showcases the use of the *ListPrf* present in *ConstrI*, which is necessary to enable us to split the list $t : \# : xs$ into t and xs . The typechecker needs some more information about t , since type families are not injective. The *split* function has type:

```

split :: ListPrf  $xs \rightarrow NP \ p \ (xs : \# : ys) \rightarrow (NP \ p \ xs, NP \ p \ ys)$ 

```

The *delCof* function is dual to *insCof*, but since we construct a *NP* indexes over $t : \# :$ xs , we need not use the *ListPrf* argument. Finally, we can assemble the application function that witnesses the semantics of *ES*:

```

applyES :: ( $\forall \ k \ . \ Eq \ (\kappa \ k) \Rightarrow ES \ \kappa \ codes \ xs \ ys \rightarrow PoA \ \kappa \ (Fix \ \kappa \ codes) \ xs$ 
            $\rightarrow Maybe \ (PoA \ \kappa \ (Fix \ \kappa \ codes) \ ys)$ 
            $\rightarrow Just \ Nil$ )
            $\rightarrow Just \ Nil$ 
applyES ES0       $\_ = Just \ Nil$ 
applyES (Ins  $\_ \ c \ es$ )  $xs = insCof \ c \ <\$> applyES \ es \ xs$ 
applyES (Del  $\_ \ c \ es$ )  $xs = delCof \ c \ xs \gg applyES \ es$ 
applyES (Cpy  $\_ \ c \ es$ )  $xs = insCof \ c \ <\$> (delCof \ c \ xs \gg applyES \ es)$ 

```

3.1.4.1 DISCUSSION

The approach of providing typed edit operations has many nice aspects. It immediately borrows the existing algorithms and metatheory and can improve the size of edit-scripts significantly by being able to provide *CpyTree*, *InsTree* and *DelTree* which copy, insert and delete entire trees instead of operating on individual constructors. This is possible because we can look at the type of the edit-script in question – substitute the insertion of a constructor by *InsTree* whenever all of its fields are also comprised solely of insertions.

Although type-safe by construction, which is undoubtedly a plus point, computing edit-scripts, with memoization, still takes $\mathcal{O}(n \times m)$ time, where n and m are the number of constructors in the source and destination trees. This means this is at least quadratic in the size of the smaller input, which is not practical for a tool that is supposed to be run multiple times per commit on large inputs (hundreds to thousands of lines). This downside is not specific to this approach, but rather quite common for tree differencing algorithms. They often belong to complexity classes that make them impractical.

Another downside comes to the surface when we want to look into merging these edit-scripts. Vassena [107] developed a merging algorithm but notes some difficult setbacks, mainly due to the heterogeneity of *ES*. Suppose, for example, we want to merge $p : ES\ xs\ ys$ and $q : ES\ xs\ zs$. This means producing an edit-script r that transforms the forest xs into some other forest ks . But how can we determine ks here? It is not always the case that there is a solution. In fact, the merge algorithm [107] for *ES* might fail due to conflicting changes or the inability to find a suitable ks . Regardless, the work of Vassena [107] was of great inspiration for this thesis in showing that there definitely is a place for type-safe approaches to differencing.

3.2 THE GENERICS-SIMPLISTIC LIBRARY

Unfortunately, the `generics-mrsop` uncovered a memory leak in the Haskell compiler itself when used for large mutually recursive families. The bugs have been reported in the GHC bug tracker² but have not been resolved at the time of writing. This means that if we wish to collect large scale real data for our experiments, we must develop an alternative approach. The `generics-simplistic` approach trades the sums-of-products structure of `generics-mrsop` for a simpler representation. Consequently, some generic functions will be more verbose than in `generics-mrsop`, but we can handle larger families without running into the aforementioned bugs.

3.2.1 THE SIMPLISTIC VIEW

The `generics-simplistic` library can be seen as a layer on top of `GHC.Generics`. Generics to ease out the definition of new generic functionality. The pattern functor approach used by `GHC.Generics`, shown in Section 2.2.1, requires the user to write a large number of typeclass instances to define even basic generic functions. Yet, the pattern functors generated by GHC are restricted to sums, products, unit, constants and metadata information. This means we can model representations as a single GADT, *SRep* defined below, indexed by the pattern functor it inhabits.

```
data SRep (φ :: * → *) :: (* → *) → * where
  S_U1 :: SRep φ U1
  S_K1 :: φ a → SRep φ (K1 i a)
  S_L1 :: SRep φ f → SRep φ (f :+: g)
  S_R1 :: SRep φ g → SRep φ (f :+: g)
  (:*) :: SRep φ f → SRep φ g → SRep φ (f :*: g)
  S_M1 :: SMeta i t → SRep φ f → SRep φ (M1 i t f)
```

²<https://gitlab.haskell.org/ghc/ghc/issues/17223> and <https://gitlab.haskell.org/ghc/ghc/issues/14987>

The handling of metadata is borrowed entirely from `GHC.Generics` and captured by the `SMeta` datatype, which records the kind of meta-information stored at the type-level.

```
data SMeta i t where
  SM_D :: Datatype d    ⇒ SMeta D d
  SM_C :: Constructor c ⇒ SMeta C c
  SM_S :: Selector s    ⇒ SMeta S s
```

The `SRep` datatype enables us to write generic functionality more concisely than `GHC.Generics`. Take the `gsize` function from Section 2.2.1 as an example. With pure `GHC.Generics`, we must use `Size` and `GSize` typeclasses. With `SRep` we can write it directly, provided we have a way to count the size of the leaves of type φ .

```
gsize :: (∀ x .  $\varphi$  x → Int) → SRep  $\varphi$  f → Int
gsize r S_UI      = 0
gsize r (S_K1 x) = r x
gsize r (S_M1 x) = gsize r x
gsize r (S_L1 x) = gsize r x
gsize r (S_R1 x) = gsize r x
gsize r (x :: y)  = gsize r x + gsize r y
```

Naturally, we still need to convert values of `GHC.Generics.Rep f x` into their closed representation, `SRep φ (GHC.Generics.Rep f)` and make some choice for φ . We could use `K1 R` as φ , essentially translating only the first layer into a generic representation, but as we shall see in Section 3.2.2, we can also translate the entire value and use a fixpoint combinator in φ .

Even though `SRep` lacks a *codes-based* approach, that is, it can be defined for arbitrary types like `GHC.Generics`, it still admits some combinators that greatly assist a programmer when writing their generic code, unlike `GHC.Generics`. The most useful are `repMap`, `repZip` and `repLeaves`, that map, zip and collect the leaves of a `SRep` respectively. These can easily be generalized to a monadic version.

```
repMap  :: (∀ x .  $\varphi$  x →  $\psi$  x) → SRep  $\varphi$  f → SRep  $\psi$  f
repZip  :: SRep  $\varphi$  f → SRep  $\psi$  f → Maybe (SRep ( $\varphi$  ::  $\psi$ ) f)
repLeaves :: SRep  $\varphi$  f → [Exists  $\varphi$ ]
```


3.2.2 MUTUAL RECURSION

The *SRep* φ *f* datatype enables us to write generic functions without resorting to type-classes and provides a simple way to interact with potentially recursive subtrees through the φ functor. Writing a deep representation involves defining a mutually recursive family as any type that is *not* a primitive type, where the choice of primitive type is parameterized through the usual κ parameter. The pseudo-code below illustrates this idea.

```
data SFix  $\kappa :: * \rightarrow *$  where
  Prim :: ( $x \in \kappa$ )  $\Rightarrow x \rightarrow$  SFix  $\kappa$  fam  $x$ 
  SFix :: ( $\neg (x \in \kappa)$ , Generic  $x$ )  $\Rightarrow$  SRep (SFix prim) (Rep  $x$ )  $\rightarrow$  SFix  $\kappa$  fam  $x$ 
```

This approach works well for simpler applications, but by defining a mutually recursive family in an *open* fashion, i.e., t is an element iff $\neg (t \in \kappa)$, for some list κ of types regarded as primitive, we would only be able to check for index equality through the *Typeable* machinery [89]. This would have to spread across the library, inherently breaking parametricity of maps and catamorphisms besides polluting the interface. Checking for index equality is crucial for the definition of many generic concepts – zippers being a prominent example, Section 3.2.4.1 – and was trivial to define in *generics-mrsop*, thanks to its *closed* approach: should two types be identified by the same index into a list containing all members of the family, then they are the same type.

To avoid having to spread *Typeables* around but still maintaining decidable type index equality we will apply the same trick here and define a family as two disjoint lists: a type-level list *fam* for the elements that belong in the family and one for the primitive types, usually denoted κ . Note that unlike *generics-mrsop*, κ here has kind $'[*]$.

Recursion is easily achieved through a *SFix* κ *fam* combinator, where *fam* $:: '[*]$ is the list of types that belong in the family and $\kappa :: '[*]$ is the list of types to be considered primitive, that is, is not unfolded into a generic representation. The *SFix* combinator has two constructors, one for carrying values of primitive types and one for unfolding a next layer of the generic representation, as defined below.

```
data SFix  $\kappa$  fam  $:: * \rightarrow *$  where
  Prim :: (PrimCnstr  $\kappa$  fam  $x$ )  $\Rightarrow x \rightarrow$  SFix  $\kappa$  fam  $x$ 
  SFix :: (CompoundCnstr  $\kappa$  fam  $x$ )  $\Rightarrow$  SRep (SFix prim) (Rep  $x$ )  $\rightarrow$  SFix  $\kappa$  fam  $x$ 
```

Here, *PrimCnstr* and *CompoundCnstr* are constraint synonyms, defined below, to encapsulate what it means for a type x to be primitive (resp. compound) with respect to the *fam* and *prim* list of types.

```
type PrimCnstr  $\kappa$  fam  $x = (Elem\ x\ \kappa, NotElem\ x\ fam)$ 
type CompoundCnstr  $\kappa$  fam  $x = (Elem\ x\ fam, NotElem\ x\ \kappa, Generic\ x)$ 
```

Elem and *NotElem* are custom constraints that state whether or not a type is an element of a list of types. They are defined with the help of the boolean type family and, in the *Elem* case, we also carry a typeclass that enables us to construct a membership proof.

```

type Elem    a as = (IsElem a as ~ 'True, HasElem a as)
type NotElem a as = IsElem a as ~ 'False

type family IsElem (a :: *) (as :: [*]) :: Bool where
  IsElem a    '[] = 'False
  IsElem a (a ' : as) = 'True
  IsElem a (b ' : as) = IsElem a as

```

HasElem a as, here, is a typeclass that produces an actual proof that the list *as* contains *a* – encoded in a datatype *ElemPrf a as*. Pattern matching on a value of type *ElemPrf a as* will unfold the structure of *as*. This is crucial in, for example, accessing typeclass instances for types in *SFix κ fam*. The *HasElem* typeclass and *ElemPrf* datatype are defined below.

```

data ElemPrf a as where
  Here  :: ElemPrf a (a ' : as)
  There :: ElemPrf a as → ElemPrf a (b ' : as)

class HasElem a as where
  hasElem :: ElemPrf a as

```

To define generic functions, we often need operation over the primitive types. We can encode this via constraints, requiring that *all* elements of κ have instances of some typeclass. Suppose we would like to write a term-level equality operator for values of type *SFix κ fam x*, as in the *Eq* typeclass. This would require to ultimately compare values of type *y*, for some *y* such that *Elem y κ*. Naturally, this can only be done if all elements of κ are members of the *Eq* typeclass. We specify that all elements of κ satisfy a constraint with the *All* [26] type family:

```

type family All c xs :: Constraint where
  All c '[]           = ()
  All c (x (' : xs)) = (c x, All c xs)

```

Now, given a function with type $(\text{All } Eq \text{ prim}) \Rightarrow \text{SFix prim } x \rightarrow \dots$, we must extract the *Eq y* instance from *All Eq prim*, for some *y* such that *IsElem y prim ~ 'True*. This is where *ElemPrf* becomes essential. By pattern matching on *ElemPrf* we are able to extract the necessary instance through the *witness* function. Naturally, once we find the instance we are looking for, we record it in a datatype for easier access. This is similar to the *Dict*

```

instance (All Eq  $\kappa$ )  $\Rightarrow$  Eq (SFix  $\kappa$  fam f) where
  (Prim x)  $\equiv$  (Prim y) = weq x y
  (SFix x)  $\equiv$  (SFix y) = maybe False (all ( $\equiv$ )  $\circ$  repLeaves) (repZip x y)

```

FIGURE 3.5: Equality instance for *SFix*.

datatype from `generics-sop`, but since we are working with a deep representation, we only use *witness* for accessing functionality over the primitive types.

```

data Witness c x where
  Witness :: (c x)  $\Rightarrow$  Witness c x
  witness ::  $\forall$  x xs c . (HasElem x xs, All c xs)  $\Rightarrow$  Proxy xs  $\rightarrow$  Witness c x
  witness _ = witnessPrf (hasElem :: ElemPrf x xs)
  where witnessPrf :: (All c xs)  $\Rightarrow$  ElemPrf x xs  $\rightarrow$  Witness c x
        witnessPrf Here      = Witness
        witnessPrf (There p) = witnessPrf p

```

The *witness* function above enables us to cast the usual (\equiv) function, from *Eq*, as operating over any element of a list of types. Pattern matching on the result of *witness* enables the compiler to access the necessary *Eq* instance. With the help of *weq* below, we define the *Eq* instance for *SFix* in Figure 3.5. Note that calling *witness* will require an explicit type annotation informing the compiler about which typeclass we wish to extract from the top-level *All* constraint.

```

weq ::  $\forall$  x xs . (All Eq xs, Elem x xs)  $\Rightarrow$  Proxy xs  $\rightarrow$  x  $\rightarrow$  x  $\rightarrow$  Bool
weq p = case witness p :: Witness Eq x of Witness  $\rightarrow$  ( $\equiv$ )

```

With the *Elem* functionality in place, we can define type-level equality for elements of a given list – given *SFix* κ fam x and *SFix* κ fam y, to be able to know whether $x \sim y$. This functionality is important when defining the zipper [45] or generic unification, and it comes for free in code-based approaches, such as `generics-mrsop`. In our current setting, we need to use the *fam* type-level list and the *HasElem* typeclass. Note that the proxies are present solely to aid the reduction of the *IsElem* type family, needed for *Elem*.

```

sameType :: (Elem x fam, Elem y fam)
           $\Rightarrow$  Proxy fam  $\rightarrow$  Proxy x  $\rightarrow$  Proxy y  $\rightarrow$  Maybe (x  $\sim$  y)
sameType _ _ _ = sameIdx (hasElem :: ElemPrf x fam) (hasElem :: ElemPrf y fam)
where sameIdx :: ElemPrf x xs  $\rightarrow$  ElemPrf x' xs  $\rightarrow$  Maybe (x  $\sim$  x')
        sameIdx Here      Here      = Just Refl
        sameIdx (There rr) (There y) = go rr y
        sameIdx _         _         = Nothing

```

```

class (CompoundCnstr  $\kappa$  fam a)  $\Rightarrow$  Deep  $\kappa$  fam a where
  dfrom :: a  $\rightarrow$  SFix  $\kappa$  fam a
  default dfrom :: (GDeep  $\kappa$  fam (Rep a))  $\Rightarrow$  a  $\rightarrow$  SFix  $\kappa$  fam a
  dfrom = SFix  $\circ$  gdfrom  $\circ$  from
  dto :: SFix  $\kappa$  fam a  $\rightarrow$  a
  default dto :: (GDeep  $\kappa$  fam (Rep a))  $\Rightarrow$  SFix  $\kappa$  fam a  $\rightarrow$  a
  dto (SFix x) = to (gdto x)

class GDeep  $\kappa$  fam f where
  gdfrom :: f x  $\rightarrow$  SRep (SFix  $\kappa$  fam) f
  gdto   :: SRep (SFix  $\kappa$  fam) f  $\rightarrow$  f x

```

FIGURE 3.6: Declaration of *Deep* and *GDeep* typeclasses.

CONVERTING TO A DEEP REPRESENTATION. Next we look at how to convert values from their shallow `GHC.Generics` representation into `generics-simplistic`. As usual, we use the *dfrom* and *dto* functions, which follow the textbook recipe of defining generic functionality with `GHC.Generics`: use a typeclass and its generic variant and use *default signatures* to bridge the gap between them. In our case, this is done with the *Deep* and *GDeep* typeclasses, declared in Figure 3.6. This technique of deriving a generic representation from some simpler generic representation is often referred to as *Generic Generic Programming* [61].

Defining the *GDeep* instances is straightforward with the exception of the *(K1 R a)* case, where we must perform different actions depending on whether or not *a* is a primitive type. Ideally we would like to write something in the lines of:

```

instance (IsElem a  $\kappa \sim$  'True)  $\Rightarrow$  GDeep  $\kappa$  fam (K1 R a) ...
instance (IsElem a  $\kappa \sim$  'False)  $\Rightarrow$  GDeep  $\kappa$  fam (K1 R a) ...

```

But GHC cannot distinguish between these two instances when resolving them. Not even `-XOverlappingInstances` can help us here. The solution, then, is to abstract the call to *IsElem* to an auxiliary typeclass, which performs a type-level pattern-match on the result of *IsElem a κ* , which we call *isPrim*, below.

```

class GDeepAtom  $\kappa$  fam (isPrim :: Bool) a where
  gdfromAtom :: Proxy isPrim  $\rightarrow$  a  $\rightarrow$  SFix  $\kappa$  fam a
  gdtoAtom   :: Proxy isPrim  $\rightarrow$  SFix  $\kappa$  fam a  $\rightarrow$  a

```

The *GDeepAtom* class possesses only two instances, one for primitive types and one for types we wish to consider as members of our mutually recursive family, which are

```

data Rose a = Fork a [Rose a]
  deriving (Eq, Show, Generic)

type PrimsRose = '[Int]
type FamRose   = '[Rose Int, [Rose Int]]

instance Deep PrimsRose FamRose (Rose Int)
instance Deep PrimsRose FamRose [Rose Int]

dfromRose :: Rose Int → SFix PrimsRose FamRose (Rose Int)
dfromRose = dfrom

dtoRose :: SFix PrimsRose FamRose (Rose Int) → Rose Int
dtoRose = dto

```

FIGURE 3.7: Usage example for *generics-simplistic*.

indicated by the *isPrim* parameter. We recall the definitions for *CompoundCnstr* and *PrimCnstr* below.

```

instance (CompoundCnstr κ fam a, Deep κ fam a) ⇒ GDeepAtom κ fam 'False a ...
instance (PrimCnstr      κ fam a)                ⇒ GDeepAtom κ fam 'True a ...

type PrimCnstr      κ fam x = (Elem x κ, NotElem x fam)
type CompoundCnstr κ fam x = (Elem x fam, NotElem x κ, Generic x)

```

Finally, the actual instance for *GDeep prim (K1 R a)* triggers the evaluation of *IsElem*, which in turn brings into scope the correct variation of the *GDeepAtom*:

```

instance (GDeepAtom κ fam (IsElem a prim) a) ⇒ GDeep κ fam (K1 R a) ...

```

With the *Deep* typeclass setup, all we have to do is declare an empty instance for every element of the family. Figure 3.7 illustrates the usage for the *Rose* datatype. The monomorphic versions of *dfrom* and *dto* simply aid the compiler by providing all necessary type parameters.

3.2.3 THE (CO)FREE (CO)MONAD

Although the *SFix* type makes for a very intuitive recursion combinator, it does not give us much flexibility: it does not support annotations nor holes. For example, suppose we want to define a generic unification algorithm: how would we represent unification variables within *SFix*? We would need an augmented *SFix* which would carry one extra constructor for unification variables. Another example would be annotating an *SFix* with some auxiliary values to make certain computations more efficient. These variations over fixpoints can be achieved by combining the free monad and the cofree

comonad in the same type, which we name *HolesAnn* κ fam φ h a . A value of type *HolesAnn* κ fam φ h a is isomorphic to a value of type a , where each constructor is annotated with φ and we might have holes of type h .

```
data HolesAnn  $\kappa$  fam  $\varphi$   $h$   $a$  where
  Hole' ::  $\varphi$   $a \rightarrow h$   $a \rightarrow$  HolesAnn  $\kappa$  fam  $\varphi$   $h$   $a$ 
  Prim'  :: (PrimCnstr  $\kappa$  fam  $a$ )       $\Rightarrow \varphi$   $a \rightarrow a \rightarrow$  HolesAnn  $\kappa$  fam  $\varphi$   $h$   $a$ 
  Roll'  :: (CompoundCnstr  $\kappa$  fam  $a$ )  $\Rightarrow \varphi$   $a \rightarrow$  SRep (HolesAnn  $\kappa$  fam  $\varphi$   $h$ ) (Rep  $a$ )
            $\rightarrow$  HolesAnn  $\kappa$  fam  $\varphi$   $h$   $a$ 
```

The *SFix* combinator presented earlier can be easily seen as the special case where annotations are the unit type, *UI*, and holes do not exist (which is captured by the empty type *VI*). We can represent all the variations over fixpoints through type synonyms:

```
type SFix       $\kappa$  fam    = HolesAnn  $\kappa$  fam UI VI
type SFixAnn    $\kappa$  fam  $\varphi$  = HolesAnn  $\kappa$  fam  $\varphi$  VI
type Holes      $\kappa$  fam    = HolesAnn  $\kappa$  fam UI
```

Again, with the help of pattern synonyms and COMPLETE pragmas – which stops GHC from issuing `-Wincomplete-patterns` warnings – we can simulate the *SFixAnn* datatype, for example.

```
pattern SFixAnn :: ()  $\Rightarrow$  (CompoundCnstr  $\kappa$  fam  $a$ )
                      $\Rightarrow \varphi$   $a \rightarrow$  SRep (SFixAnn  $\kappa$  fam  $\varphi$ ) (Rep  $a$ )  $\rightarrow$  SFixAnn  $\kappa$  fam  $\varphi$   $a$ 
pattern SFixAnn ann  $x$  = Roll' ann  $x$ 
pattern PrimAnn  :: ()  $\Rightarrow$  (PrimCnstr  $\kappa$  fam  $a$ )  $\Rightarrow \varphi$   $a \rightarrow a \rightarrow$  SFixAnn  $\kappa$  fam ann  $a$ 
pattern PrimAnn ann  $x$  = Prim' ann  $x$ 
{-# COMPLETE SFixAnn , PrimAnn #-}
```

3.2.3.1 ANNOTATED FIXPOINTS

Catamorphisms are used in a large number of computations over recursive structures. They receive an algebra that is used to consume one layer of a datatype at a time and consumes the whole value of the datatype using this *recipe*. The definition of the catamorphism is trivial in a setting where we have explicit recursion:

```
cata :: ( $\forall b$  . (CompoundCnstr  $\kappa$  fam  $b$ )  $\Rightarrow$  SRep  $\varphi$  (Rep  $b$ )  $\rightarrow \varphi$   $b$ )
        $\rightarrow$  ( $\forall b$  . (PrimCnstr  $\kappa$  fam  $b$ )  $\Rightarrow b \rightarrow \varphi$   $b$ )  $\rightarrow$  SFix  $\kappa$  fam  $h$   $a \rightarrow \varphi$   $a$ 
cata  $f$   $g$  (SFix  $x$ ) =  $f$  (repMap (cata  $f$   $g$ )  $x$ )
cata  $_$   $g$  (Prim  $x$ ) =  $g$   $x$ 
```

One example of catamorphisms is computing the *height* of a recursive structure. It can be defined with *cata* in a simple manner with the help of the *Const* functor.

```
newtype Const t x = Const {getConst :: t}
heightAlgebra :: SRep (Const Int) xs → Const Int iy
heightAlgebra = Const ∘ (1+) ∘ maximum ∘ (0:) ∘ map (exElim getConst) ∘ repLeaves
height :: SFix κ fam a → Int
height = getConst ∘ cata heightAlgebra
```

Now imagine our particular application makes a number of decisions based on the height of the (generic) trees it handles. Calling *height* at each of those decision points is time consuming. It is much better to compute the height of a tree only once and keep the intermediary results annotated in their respective subtrees. We can easily do so with our *SFixAnn* *cofree comonad* [37]. In fact, we would say that the height is a synthesized attribute in *attribute grammar* [52] lingo.

```
synthesize :: (∀ b . (CompoundCnstr κ fam a) ⇒ SRep φ (Rep b) → φ b)
           → (∀ b . (PrimCnstr κ fam a) ⇒ b → φ b)
           → SFix κ fam a → SFixAnn κ fam φ a
synthesize f g = cata (λr → SFixAnn (f (repMap getAnn r)) r) (λa → PrimAnn (g b) b)
```

Finally, an algorithm that constantly queries the height of the subtrees can be computed in two passes: in the first pass we compute the heights and leave them annotated in the tree, in the second we run the algorithm. Moreover, we can compute all the necessary synthesized attributes an algorithm needs in a single preprocessing phase. This is a crucial maneuver to make sure our generic programs can scale to real-world inputs. Naturally, *cata* and *synthesize* are actually implemented in their monadic form and over *HolesAnn* for maximal generality.

3.2.4 PRACTICAL FEATURES

Whilst developing *hdiff* (Chapter 5), we ran into a number of practicalities regarding the underlying generic programming library. Of particular importance are zippers and unification, which play a big role in the algorithms underlying the *hdiff* approach. This section gives an overview of those features.

3.2.4.1 ZIPPERS

Zippers [45] are a well established technique for traversing a recursive data structure keeping track of a focus point. Defining generic zippers is not new, this has been done by many authors [1, 43, 112] for many different classes of datatypes in the past. In our

particular case, we are not interested in traversing a generic representation by means of the usual zipper traversals – up, down, left and right – which move the focus point. Instead, we just want a datatype that encodes a context with one focus, encoded by *SZip* below. A value of type *SZip ty w f* represents a value of type *SRep w f* with one *hole*, or *focus*, in a position with type *ty*.

```
data SZip ty w f where
  Z_L1  :: SZip ty          wf → SZip ty w (f :+ : g)
  Z_R1  :: SZip ty          wg → SZip ty w (f :+ : g)
  Z_PairL :: SZip ty w f → SRep w g → SZip ty w (f :* : g)
  Z_PairR :: SRep w f → SZip ty w g → SZip ty w (f :* : g)
  Z_M1   :: SMeta i t → SZip ty w f → SZip ty w (M1 i t f)
  Z_KH   ::              → SZip ty w (K1 i a)
```

The *Zipper* datatype will ensure that the focus lies in a recursive position. Its definition is given below. It encapsulates the *ty* above as an existential type and keeps the focus point accessible. We also pass around a constraint-kinded variable to enable one to specify custom constraints about the types in question.

```
data Zipper c f g t where
  Zipper :: c ⇒ SZip t f (Rep t) → g t → Zipper c f g t
```

Given a value of type *SZip ty φ t* and a value of type *φ ty*, it is straightforward to plug the hole and produce a *SRep φ t*. The other way around, however, is more complicated. Given a *SRep φ t*, we might have many possible zippers – binary trees, for example, can have a hole on the left or on the right branch. Consequently, we must return a list of zippers. The *zippers* function below does exactly that. Its type is convoluted because it works over *HolesAnn* (and therefore also for *SFix*, *SFixAnn* and *Holes*), but it is conceptually simple: given a test for whether a hole of type *φ a* is actually a hole in a recursive position, we return the list of possible zippers for a value with holes. The definition is standard and we encourage the interested reader to check the source code for more details (Appendix A).

```
type Zipper' κ fam φ h t
  = Zipper (CompoundCnstr κ fam t) (HolesPhi κ fam φ h) (HolesPhi κ fam φ h) t
zippers :: (∀ a . (Elem t fam) ⇒ h a → Maybe (a :~ : t))
  → HolesPhi κ fam φ h t → [Zipper' κ fam φ h t]
```


3.2.4.2 UNIFICATION AND ANTI-UNIFICATION

Both unification and anti-unification algorithms make up an important part of the vernacular of term-manipulation. Unsurprisingly, we will also need to implement these features into `generics-simplistic`. We use them extensively in Chapter 5. This section provides an overview of the (anti-)unification provided by `generics-simplistic`.

Syntactic unification algorithms [94] receive as input two terms t and u with variables and outputs substitutions σ such that $\sigma t \equiv \sigma u$, when such σ exists. Anti-unification[91], on the other hand, receives two terms t and u and outputs one term r and two substitutions σ and φ such that $t \equiv \sigma r$ and $u = \varphi r$.

With our current setup, we want to unify two terms of type `Holes κ fam φ at`, that is, two elements of the mutually recursive family `fam` with unification variables of type φ . A substitution is given by:

```
type Subst  $\kappa$  fam  $\varphi$  = Map (Exists  $\varphi$ ) (Exists (Holes  $\kappa$  fam  $\varphi$ ))
```

We need the existentials here in order to use the builtin, homogeneous, `Data.Map`. Naturally, when looking for the value associated with a key within the substitution we will run into a type error as soon as we unwrap the `Exists`. There are a number of solutions to this. For one, we could use the `sameTy` function and ensure they are of the same type. Pragmatically though, as long as we ensure we only insert keys φ at associated with values `Holes κ fam φ at`, the type indexes will never differ and we can safely call `unsafeCoerce` to mitigate any performance overhead. We chose to use `unsafeCoerce` but stress that it can be easily avoided with a call to `sameTy`.

```
substInsert :: (Ord (Exists  $\varphi$ )) => Subst  $\kappa$  fam  $\varphi$  ->  $\varphi$  at -> Holes  $\kappa$  fam  $\varphi$  at
              -> Subst  $\kappa$  fam  $\varphi$ 

substLkup  :: (Ord (Exists  $\varphi$ )) => Subst  $\kappa$  fam  $\varphi$  ->  $\varphi$  at -> Maybe (Holes  $\kappa$  fam  $\varphi$  at)
```

When attempting to solve a unification problem, there are two types of failures that can occur: symbol clashes happen when we try to unify different symbols, for example, $c\ x$ is not unifiable with $c'\ x$ because $c \not\equiv c'$; and occurs-check errors are raised when there is a loop in the substitution. For example, if we try to unify $c\ (c'\ x)$ with $c\ x$, we would have to substitute x for $c'\ x$, but this would never terminate. We encode these errors in the `UnifyErr` datatype, making it easy for users of the library to catch these errors and extract information from them.

```
data UnifyErr  $\kappa$  fam  $\varphi$  where
  OccursCheck :: [Exists  $\varphi$ ] -> UnifyErr  $\kappa$  fam  $\varphi$ 
  SymbolClash :: Holes  $\kappa$  fam  $\varphi$  at -> Holes  $\kappa$  fam  $\varphi$  at -> UnifyErr  $\kappa$  fam  $\varphi$ 
```

```

lgg :: (All Eq κ) ⇒ Holes κ fam φ at → Holes κ fam ψ at
      → Holes κ fam (Holes κ fam φ :*: Holes κ fam ψ) at
lgg (Prim x) (Prim y) =
  | weq (Proxy :: Proxy κ) x y = Prim x
  | otherwise                    = (Prim x :*: Prim y)
lgg x@(Roll rx) y@(Roll ry) = case zipSRep rx ry of
  Nothing → Hole (x :*: y)
  Just r  → Roll (repMap (uncurry' lgg) r)
lgg x y = Hole (x :*: y)

```

FIGURE 3.8: Classic anti-unification algorithm [91], producing the least general generalization of two trees.

The *unify* function has the type one would expect: given two terms with unification variables of type ϕ , either they are not unifiable or there exists a substitution that makes them equal.

```

unify :: (Ord (Exists φ), EqHO φ) ⇒ Holes κ fam φ at → Holes κ fam φ at
      → Except (UnifyErr κ fam φ) (Subst κ fam φ)

```

Our *unify* function is a constraint-based unifier which computes the most general unifier in two phases: first it collects all the necessary equivalences, then it tries to produce an idempotent substitution from the gathered equivalences. We omit technical details regarding the implementation of the unification algorithm and refer the reader to the existing literature [94].

Anti-unification [91] is dual to unification. It is the process of identifying the longest common prefix of two terms. For example, take $x = \text{Bin } (\text{Bin } 1\ 2)\ \text{Leaf}$ and $y = \text{Bin } (\text{Bin } 1\ 3)\ (\text{Bin } 4\ 5)$, the term $\text{Bin } (\text{Bin } 1\ a)\ b$ is the least general generalization of x and y . That is, there exist two instantiations of a and b yielding x or y . The term $\text{Bin } c\ b$ is also a generalization of x and y , but it is not the *least* general because to obtain x or y we would have instantiate c as $\text{Bin } 1\ 2$ or $\text{Bin } 1\ 3$, and these terms can be further anti-unified. Figure 3.8 illustrates the implementation of the syntactical anti-unification algorithm.

3.3 DISCUSSION

In this chapter we explored two different ways of writing generic programs that must work over mutually recursive families. Looking back at the spectrum of generic programming libraries, in Figure 2.11, we had an open hole for *code-based* approach with ex-

	Pattern Functors	Codes
No Explicit Recursion	<code>GHC.Generics</code>	<code>generics-sop</code>
Simple Recursion	<code>generics-simplistic</code>	<code>generics-mrsop</code>
Mutual Recursion		

FIGURE 3.9: *Updated spectrum of generic programming libraries.*

plicit recursion of any type, which can be filled by `generics-mrsop`. When it comes to pattern functors, although `regular` and `multirec` already exist, using those libraries imposes a significant overhead when compared to `generics-simplistic`, for they do not support combinator-based generic programming. The updated table of generic programming libraries is given in Figure 3.9, where we place our libraries in the spectrum of generic programming variants.

Unfortunately, the `generics-mrsop` heavy usage of type families triggers a memory leak in the compiler. This renders the library unusable for large families of mutually recursive datatypes at the time of writing this thesis. Luckily, however, we were able to work around that by dropping the sums of products structure but maintaining a combinator-based approach in `generics-simplistic`, which enabled us to run our experiments with real-world data, as discussed in Chapter 6.

Although both `generics-mrsop` and `generics-simplistic` are quite similar, the sums-of-products structure of `generics-mrsop` is more convenient to program with. We must pattern match on less cases (first on sums, then on products) and we have access to more combinators to dissect the sums-of-products structure. Now, that being said, not all generic programming applications will require that level of control and some could still get by without it. The `hdiff` (Chapter 5) is an example of the later. Our first implementation with `generics-mrsop` was more natural, but we were still able to implement it using `generics-simplistic`. The advantages of `generics-simplistic` are twofold: (A) it uses significantly less type families and does not run into the bugs that `generics-mrsop` did and (B) it is simpler to use and it piggybacks on a number of familiar features from `GHC.Generics`.

While developing the `generics-mrsop` and `generics-simplistic` libraries, which happened under close collaboration with Alejandro Serrano, we also explored a number of variants of these libraries such as `kind-generics` [98], which enables a user to represent almost any Haskell datatype generically, including *GADTs*. These are out of the scope of this thesis since we do not require all of that expressivity to write our differencing algorithms.



STRUCTURAL PATCHES

The `gdiff` [55] approach, discussed in Section 3.1.4, which flattens a tree into a list, following classical tree edit distance algorithms encoded through using type-safe edit-scripts, inherits the problems of edit-script based approaches. These include ambiguity on the representation of patches, non-uniqueness of optimal solutions and difficulty of merging. The `stdiff` approach, discussed through this chapter, arose from our study of the difficulties about merging `gdiff` patches [107].

The heterogeneity of Patch_{GD} makes merging difficult. Recall that a value of type $\text{Patch}_{\text{GD}}\ x\ y$ transforms a list of trees x s into a list of trees y s. If we are given two patches $\text{Patch}_{\text{GD}}\ x\ y$ and $\text{Patch}_{\text{GD}}\ x\ z$, we would like to produce two patches $\text{Patch}_{\text{GD}}\ y\ r$ and $\text{Patch}_{\text{GD}}\ z\ r$ such that the canonical merge square commutes. The problem becomes clear when we try to determine r s correctly: sometimes such r s might not even exist [107].

Our `stdiff` approach, or, *structural patches*, marks our first attempt at defining a *type-indexed* patch datatype, Patch_{ST} , in pursuit of better behaved merge algorithms. The overall idea consists in making sure that the type of patches is also *tree structured*, as opposed to managing a list-like patch data structure that is supposed to operate over tree structured data. As it turns out, it is not possible to have fully homogeneous patches, but we were able to identify homogeneous parts of our patches which we can use to synchronize changes when defining our merge operation, but let us not get ahead of ourselves.

Structural Patches differ from edit-scripts by using tree-shaped, homogeneous patch – a patch transforms two values of the same type. The edit operations themselves are analogous to edit scripts, we support insertions, deletions and copies, but these are structured to follow the sums-of-products of datatypes: there is one way of changing sums,

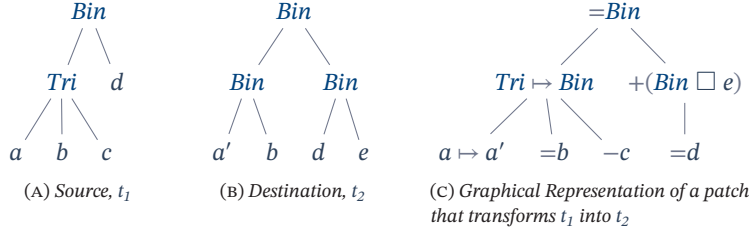


FIGURE 4.1: Graphical representation of a simple transformation. Copies, insertions and deletions around the tree are represented with $=$, $+$ and $-$ respectively. Modifications are denoted $\cdot \mapsto \cdot$.

one way of changing products and one way of changing the recursive positions of a value. For example, consider the following trees:

$$t_1 = \text{Bin} (\text{Tri } a \ b \ c) \ d$$

$$t_2 = \text{Bin} (\text{Bin } a' \ b) \ (\text{Bin } d \ e)$$

These are the trees that are depicted in Figure 4.1(A) and Figure 4.1(B) respectively. How should we represent the transformation mapping t_1 into t_2 ? Traversing the trees from their roots, we see that on the outermost level they both consist of a *Bin*, yet the fields of the source and destination nodes are different: the first field changes from a *Bin* to a *Tri*, which requires us to reconcile the list of fields $[a, b, c]$ into $[a', b]$. This can be done by the means of an edit-script. The second field, however, witnesses a change in the recursive structure of the type. We see that we have inserted new information, namely $(\text{Bin} \ \square \ e)$. After inserting this *context*, we simply copy d from the source to the destination. This transformation has been sketched graphically in Figure 4.1(C), and showcases all the necessary pieces we will need to write a general encoding of transformations between objects that support insertions, deletions and copies.

The stdiff approach to differencing is unlike the edit-scripts we saw previously, using the shape of the datatype in question to define a structured notion of patch. As we will see in the remainder of this chapter, however, *computing* these patches is intractable as the number of possible patches explodes. This lead us to abandon this approach in favor of the differencing algorithm presented in Chapter 5. Nonetheless, we believe there is value in studying this approach. For one thing it explores a different part in the design space compared to the *gdif* algorithm we saw previously, but it also provides insights that help understand the more efficient approach in Chapter 5.

To write the *stdiff* algorithms in Haskell, we must rely on the *generics-mrsop* library (Section 3.1) as our generic programming workhorse for two reasons. First, we do require the concept of explicit sums of products in the very definition of $\text{Patch}_{\text{ST}} \ x$.

Secondly, we need `gdiff`'s assistance in computing patches (Section 4.3.2) and `gdiff` also requires, to a lesser extent, sums of products structured datatypes, hence is easily written with `generics-mrsop`, as seen in Section 3.1.4.

The contributions in this chapter arise from joint work with Pierre-Evariste Dagand, published in TyDe 2017 [74] and coded in Agda [83]repository¹. Later, we collaborated closely with Arian van Putten, in translating the Agda code to Haskell, extending its scope to mutually recursive datatypes. The code presented here, however, is loosely based on Van Putten's translation of our Agda repository to Haskell as part of his Master thesis work [92]. We chose to present all of our work in a single programming language to keep the thesis consistent throughout.

In this chapter we will delve into the construction of `PatchST` and its respective components. Firstly, we familiarize ourselves with `PatchST` and its application function (Section 4.1). Next we look into merging and its commutativity proof in Section 4.2. Lastly, we discuss the `diff` function in Section 4.3, which comprises a significant drawback of the `stdiff` approach for its computational complexity.

4.1 THE TYPE OF PATCHES

Next we look at the `PatchST` type, starting with a single layer of datatype, i.e., a single application of the datatypes pattern functor. Later, in Section 4.1.2 we extend this treatment to recursive datatypes, essentially by taking the fixpoint of the constructions in Section 4.1.1. The `generics-mrsop` library (Chapter 3) will be used throughout the exposition.

Recall that a datatype, when seen through its initial algebra semantics [105], can be seen as an infinite succession of applications of its pattern functor, call it F , to itself: $\mu F = F(\mu F)$. The `PatchST` type will describe the differences between values of μF by successively applying the description of differences between values of type F , closely following the initial algebra semantics of datatypes.

4.1.1 FUNCTORIAL PATCHES

Handling *one layer* of recursion is done by addressing the possible changes at the sum level, followed by some reconciliation at the product level when needed.

The first part of our algorithm handles the *sums* of the universe. Given two values, x and y , it computes the *spine*, capturing the largest common coproduct structure. We distinguish three possible cases:

¹<https://github.com/VictorCMiraldo/stdiff>

- x and y are fully equal, in which case we copy the full values regardless of their contents. They must also be of the same type.
- x and y have the same constructor – i.e., $x = \text{inj } c \text{ } px$ and $y = \text{inj } c \text{ } py$ – but some subtrees of x and y are distinct, in which case we copy the head constructor and handle all arguments pairwise.
- x and y have distinct constructors, in which case we record a change in constructor and a choice of the alignment of the source and destination’s constructor fields. Here, x and y might be of a different type in the family.

The datatype *Spine*, defined below, formalizes this description. The three cases we describe above correspond to the three constructors of *Spine*. When two values are not equal, we need to represent the differences somehow. If the values have the same constructor we need to reconcile the fields of that constructor whereas if the values have different constructors we need to reconcile the products that make the fields of the constructors. We index the datatype *Spine* by the sum codes it operates over because we need to lookup the fields of the constructors that have changed, and *align* them in the case of *SChg*. Alignments will be introduced shortly. For the time being, let us continue to focus on spines. Intuitively, spines act on sums and capture the “largest shared co-product structure”. Recall $\kappa :: \text{kon} \rightarrow *$ interprets the opaque types in the mutually recursive family in question and $\text{codes} :: [[\text{Atom kon}]]$ lists all the sums-of-products in the family. Both come from *generics-mrsop* representation of mutually recursive datatypes, discussed in Section 3.1.

```

data Spine  $\kappa$   $\text{codes} :: [[\text{Atom kon}]] \rightarrow [[\text{Atom kon}]] \rightarrow *$  where
  Scp  :: Spine  $\kappa$   $\text{codes}$   $s_1$   $s_1$ 
  SCns :: Constr  $s_1$   $c_1 \rightarrow \text{NP } (\text{At } \kappa \text{ codes}) (\text{Lkup } c_1 \text{ } s_1) \rightarrow \text{Spine } \kappa \text{ codes } s_1 \text{ } s_1$ 
  SChg :: Constr  $s_1$   $c_1 \rightarrow \text{Constr } s_2 \text{ } c_2 \rightarrow \text{Al } \kappa \text{ codes } (\text{Lkup } c_1 \text{ } s_1) (\text{Lkup } c_2 \text{ } s_2)$ 
          $\rightarrow \text{Spine } \kappa \text{ codes } s_1 \text{ } s_2$ 

```

Our Agda model [74] handles only regular types, or, mutually recursive families consisting of a single datatype. Hence, the *Spine* type would arise naturally as a homogeneous type. While extending the Agda model to a full fledged Haskell implementation, together with Van Putten [92], we noted how this would severely limit the number of potential copy opportunities throughout patches. For example, imagine we want to patch the following values:

```

data  $T = T_1 \text{ } XYZ | T_2 \text{ } U$ 
data  $U = U_1 \text{ } XYZ | U_2 \text{ } T$ 
 $\text{diff } (T_1 \text{ } x_1 \text{ } y_1 \text{ } z_1) (U_1 \text{ } x_2 \text{ } y_2 \text{ } z_2) = \text{SChg } T_1 \text{ } U_1 \text{ } \dots$ 

```

With a fully homogeneous *Spine* type, our only option is to delete T_1 , then insert U_1 at the *recursion* layer (4.1.2) This would be unsatisfactory as it only allows copying

of one of the fields, where `gdifff` would be able to copy more fields for it does not care about the recursive structure.

The semantics of *Spine* are straightforward, but before continuing with *applySpine*, a short technical interlude is necessary. The *testEquality*, below, is used to compare the type indices for propositional equality. It comes from *Data.Type.Equality* and has type $f\ a \rightarrow f\ b \rightarrow \text{Maybe}\ (a \sim\!:\! b)$. Also note that we must pass two *SNat* arguments to disambiguate the *ix* and *iy* type variables. Without those arguments, these variables would only appear as an argument to a type family, which may not be injective and would trigger a type error. Using the *SNat* singleton [29] is the standard Haskell type-level programming workaround to this problem.

```
data SNat :: Nat → * where ...
```

The *applySpine* function is given by checking the provided value is made up with the required constructor. In the *SCns* case we must ensure that type indices match – for Haskell type families may not be injective – then simply map over the fields with the *applyAt* function, which applies changes to atoms. Otherwise, we reconcile the fields with the *applyAl* function, whose definition follows shortly.

```
applySpine :: (EqHO κ) ⇒ SNat ix → SNat iy
  → Spine κ codes (Lkup ix codes) (Lkup iy codes)
  → Rep κ (Fix κ codes) (Lkup ix codes)
  → Maybe (Rep κ (Fix κ codes) (Lkup iy codes))

applySpine _ _ Scp x = return x
applySpine ix iy (SCns c1 dxs) (sop → Tag c2 xs) = do
  Refl ← testEquality ix iy
  Refl ← testEquality c1 c2
  inj c2 <$> (mapNPM applyAt (zipNP dxs xs))
applySpine _ _ (SCHg c1 c2 al) (sop → Tag c3 xs) = do
  Refl ← testEquality' c1 c3
  inj c2 <$> applyAl al xs
```

The *Spine* datatype and *applySpine* are responsible for matching the *constructors* of two trees, but we still need to determine how to continue representing the difference in the products of data stored therein. At this stage in our construction, we are given two heterogeneous lists, corresponding to the fields associated with two distinct constructors. As a result, these lists need not have the same length nor store values of the same type. To do so, we need to decide how to line up the constructor fields of the source and destination. We shall refer to the process of reconciling the lists of constructor fields as solving an *alignment* problem.

Finding a suitable alignment between two lists of constructor fields amounts to finding a suitable *edit-script*, that relates source fields to destination fields. The *Al* datatype

below describes such edit-scripts for a heterogeneously typed list of atoms. These scripts may insert fields in the destination (*Ains*), delete fields from the source (*Adel*), or associate two fields from both lists (*AX*).

```

data Al  $\kappa$  codes :: [Atom kon] → [Atom kon] → * where
  A0   :: Al  $\kappa$  codes '[]' '[]'
  AX   :: At  $\kappa$  codes x           → Al  $\kappa$  codes xs ys → Al  $\kappa$  codes (x '': xs) (x '': ys)
  Adel :: NA  $\kappa$  (Fix  $\kappa$  codes) x → Al  $\kappa$  codes xs ys → Al  $\kappa$  codes (x '': xs)      ys
  Ains :: NA  $\kappa$  (Fix  $\kappa$  codes) x → Al  $\kappa$  codes xs ys → Al  $\kappa$  codes      xs (x '': ys)

```

We require alignments to preserve the order of the arguments of each constructor, thus forbidding permutations of arguments. In effect, the datatype of alignments can be viewed as an intentional representation of (partial) *order and type preserving maps*, along the lines of McBride’s order preserving embeddings [65], mapping source fields to destination fields. This makes sure that our patches also give rise to tree mappings (Section 2.1.2) in the classical tree-edit distance sense.

Given a partial embedding for atoms, we can therefore interpret alignments into a function transporting the source fields over to the corresponding destination fields, failure potentially occurring when trying to associate incompatible atoms. Recall (\times) and ϵ are the constructors of type *NP* (Page 23):

```

applyAl :: (EqHO  $\kappa$ ) ⇒ Al  $\kappa$  codes xs ys → PoA  $\kappa$  (Fix  $\kappa$  codes) xs
           → Maybe (PoA  $\kappa$  (Fix  $\kappa$  codes) ys)
applyAl A0            $\epsilon$            = return  $\epsilon$ 
applyAl (AX dx dxs) (x  $\times$  xs) = ( $\times$ ) <$> applyAt (dx :: x) <*> applyAl dxs xs
applyAl (Ains x dxs)      xs = (x  $\times$ ) <$> applyAl dxs xs
applyAl (Adel x dxs) (y  $\times$  xs) = guard (eq1 x y) <*> applyAl dxs xs

```

Finally, when synchronizing atoms we must distinguish between a recursive position or opaque data. In case of opaque data, we simply record the old value and the new value.

```

data TrivialK ( $\kappa$  :: kon → *) :: kon → * where
  Trivial ::  $\kappa$  kon →  $\kappa$  kon → TrivialK  $\kappa$  kon

```

In case we are at a recursive position, we record a potential change in the recursive position with *Al μ* , which we will get to shortly.

```

data At ( $\kappa$  :: kon → *) (codes :: [[Atom kon]]) :: Atom kon → * where
  AtSet  :: TrivialK  $\kappa$  kon → At  $\kappa$  codes ('K kon)
  AtFix  :: (IsNat ix) ⇒ Al $\mu$   $\kappa$  codes ix ix → At  $\kappa$  codes ('I ix)

```

The application function for atoms follows the same structure. In case we are applying a patch to an opaque type, we must understand whether said patch represents a copy, i.e., the source and destination values are the same. If that is the case, we simply copy the provided value. Otherwise, we must ensure the provided value matches the source value. The recursive position case is directly handled by the *applyAlμ* function.

```

applyAt :: (EqHO ki) ⇒ At κ codes at → NA κ (Fix κ codes) at
          → Maybe (NA κ (Fix κ codes) at)
applyAt (AtSet (Trivial x y)) (NAK a)
  | eqHO x y = Just (NAK a)
  | eqHO x a = Just (NAK b)
  | otherwise = Nothing
applyAt (AtFix px) (NAI x) = NAI <$> applyAlμ px x

```

The last step is to address how to make changes over the recursive structure of our value, defining *Alμ* and *applyAlμ*, which will be our next concern.

4.1.2 RECURSIVE CHANGES

In the previous section, we presented patches describing changes to the coproducts, products and atoms of our *SoP* universe. This treatment handled just a single layer of the fix-point construction. In this section, we tie the knot and define patches describing changes to arbitrary *recursive* datatypes.

To represent generic patches on values of *Fix codes ix*, we will define two mutually recursive datatypes *Alμ* and *Ctx*. The semantics of both these datatypes will be given by defining how to *apply* them to arbitrary values:

- Much like alignments for products, a similar phenomenon appears at fixpoints. When comparing two recursive structures, we can insert, remove or modify constructors. Since we are working over mutually recursive families, removing or inserting constructors can change the overall type. We will use *Alμ ix iy* to specify these edit-scripts at the constructor-level, describing a transformation from *Fix codes ix* to *Fix codes iy*.
- Whenever we choose to insert or delete a recursive subtree, we must specify *where* this modification takes place. To do so, we will define a new type *Ctx ...* :: '[*Atom kon*] → *, inspired by zippers [45, 64], to navigate through our data-structures. A value of type *Ctx ... p* selects a single atom *I* from the product of type *p*.

Modeling changes over fixpoints closely follows our definition of alignments of products. Instead of inserting and deleting elements of the product we insert, delete or modify *constructors*. Our previous definition of spines merely matched the constructors of the source and destination values – but never introduced or removed them. It is precisely these operations that we must account for here.

```

data  $Al\mu \kappa \text{ codes} :: Nat \rightarrow Nat \rightarrow *$  where
   $Spn :: Spine \kappa \text{ codes} (Lkup \text{ ix codes}) (Lkup \text{ iy codes})$ 
     $\rightarrow Al\mu \kappa \text{ codes ix iy}$ 
   $Ins :: Constr (Lkup \text{ iy codes}) c \rightarrow InsCtx \kappa \text{ codes ix} (Lkup \text{ c} (Lkup \text{ iy codes}))$ 
     $\rightarrow Al\mu \kappa \text{ codes ix iy}$ 
   $Del :: Constr (Lkup \text{ ix codes}) c \rightarrow DelCtx \kappa \text{ codes iy} (Lkup \text{ c} (Lkup \text{ ix codes}))$ 
     $\rightarrow Al\mu \kappa \text{ codes ix iy}$ 

```

The first constructor, *Spn*, does not perform any new insertions and deletions, but instead records a spine and an alignment of the underlying product structure. This closely follows the patches we have seen in the previous section. To insert a new constructor, *Ins*, requires two pieces of information: a choice of the new constructor to be introduced, called *c*, and the fields associated with that constructor. Note that we only need to record *all but one* of the constructor’s fields, as represented by a value of type $InsCtx \text{ ki codes ix} (Lkup \text{ c} (Lkup \text{ iy codes}))$. Deleting a constructor is analogous to insertions, with *InsCtx* and *DelCtx* being slight variations over *Ctx*, where one actually flips the arguments to ensure the transformation is on the right direction.

```

type  $InsCtx \kappa \text{ codes} = Ctx \kappa \text{ codes} \quad (Al\mu \kappa \text{ codes})$ 
type  $DelCtx \kappa \text{ codes} = Ctx \kappa \text{ codes} (Flip (Al\mu \kappa \text{ codes}))$ 
newtype  $Flip f \text{ ix iy} = Flip \{unFlip :: f \text{ iy ix}\}$ 

```

Our definition of insertion and deletions relies on identifying *one* recursive argument among the product of possibilities. To model this accurately, we define an indexed zipper to identify a recursive atom (indicated by a value of type *I*) within a product of atoms. Conversely, upon deleting a constructor from the source structure, we exploit *Ctx* to indicate the subtree that should be used for the remainder of the patch application, discarding all other constructor fields. We parameterize the *Ctx* type with a $Nat \rightarrow Nat \rightarrow *$ argument to distinguish between these two cases, as seen above.

```

data  $Ctx \kappa \text{ codes} (p :: Nat \rightarrow Nat \rightarrow *) (ix :: Nat) :: [Atom \text{ kon}] \rightarrow *$  where
   $H :: (IsNat \text{ iy}) \Rightarrow p \text{ ix iy} \rightarrow PoA \kappa (Fix \kappa \text{ codes}) \text{ xs} \rightarrow Ctx \kappa \text{ codes p ix} ('I \text{ iy} \text{ '}: \text{ xs})$ 
   $T :: NA \kappa (Fix \kappa \text{ codes}) a \rightarrow Ctx \kappa \text{ codes p ix xs} \rightarrow Ctx \kappa \text{ codes p ix} (a \text{ '}: \text{ xs})$ 

```

Consequently, we will have two application functions for contexts, one that inserts and one that removes contexts. This makes the need to flip the type indexes of *Al μ* , when

defining *DelCtx*, clear. Inserting a context is done by receiving a tree and returning the product stored in the context with the distinguished field applied to the received tree:

$$\begin{aligned}
\text{insCtx} &:: (\text{IsNat } ix, \text{EqHO } \kappa) \Rightarrow \text{InsCtx } \kappa \text{ codes } ix \text{ xs} \rightarrow \text{Fix } \kappa \text{ codes } ix \\
&\rightarrow \text{Maybe } (\text{PoA } \kappa (\text{Fix } \kappa \text{ codes } ix)) \\
\text{insCtx } (H \ x \ rest) \ v &= (\times_{rest}) \circ \text{NA}_I \langle \$ \rangle \text{ applyAl}\mu \ x \ v \\
\text{insCtx } (T \ a \ ctx) \ v &= (a \times) \langle \$ \rangle \text{ insCtx } ctx \ v
\end{aligned}$$

The deletion function discards any information we have about all the constructor fields, except for the subtree used to continue the patch application process. This is a consequence of our design decision, at the time, of having application functions as permissive as possible. Intuitively, the deletion context identifies the only field that should not be deleted. By not checking whether the elements we are applying to match the ones that should be deleted, we get an application function that applies to more elements for free.

$$\begin{aligned}
\text{delCtx} &:: (\text{IsNat } ix, \text{EqHO } \kappa) \Rightarrow \text{DelCtx } \kappa \text{ codes } ix \text{ xs} \rightarrow \text{PoA } \kappa (\text{Fix } \kappa \text{ codes } ix) \\
&\rightarrow \text{Maybe } (\text{Fix } \kappa \text{ codes } ix) \\
\text{delCtx } (H \ x \ rest) \ (\text{NA}_I \ v \times p) &= \text{applyAl}\mu \ (\text{unFlip } x) \ v \\
\text{delCtx } (T \ a \ ctx) \ (at \ \times p) &= \text{delCtx } ctx \ p
\end{aligned}$$

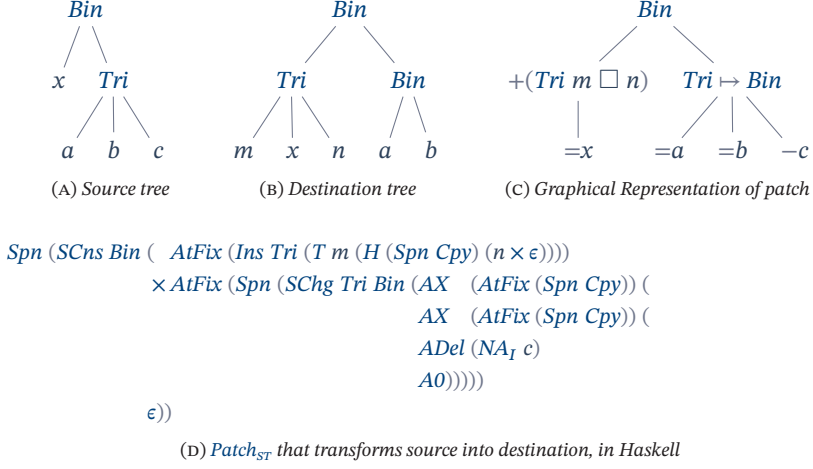
Finally, the application function for *Al* μ is nothing but selecting whether we should use the spine functionality or insertion and deletion of a context.

$$\begin{aligned}
\text{applyAl}\mu &:: (\text{IsNat } ix, \text{IsNat } iy, \text{EqHO } \kappa) \Rightarrow \text{Al}\mu \ \kappa \text{ codes } ix \text{ iy} \rightarrow \text{Fix } \kappa \text{ codes } ix \\
&\rightarrow \text{Maybe } (\text{Fix } \kappa \text{ codes } iy) \\
\text{applyAl}\mu \ (\text{Spn } sp) \ (\text{Fix } rep) &= \text{Fix } \langle \$ \rangle \text{ applySpine } _ _ \text{ spine } rep \\
\text{applyAl}\mu \ (\text{Ins } c \ ctx) \ (\text{Fix } rep) &= \text{Fix } \circ \text{inj } c \ \langle \$ \rangle \text{ insCtx } ctx \ f \\
\text{applyAl}\mu \ (\text{Del } c \ ctx) \ (\text{Fix } rep) &= \text{delCtx } ctx \ \langle \$ \rangle \text{ match } c \ rep
\end{aligned}$$

The two underscores at the *Spn* case are just an extraction of the necessary singletons to make the *applySpine* typecheck. These can be easily replaced by *getSNat* with the correct proxies. Figure 4.2 provides a graphical illustration of a value of type *Patch_{ST}* that transforms two concrete trees.

$$\begin{aligned}
\text{type } \text{Patch}_{\text{ST}} \ \kappa \text{ codes } ix &= \text{Al}\mu \ \kappa \text{ codes } ix \text{ ix} \\
\text{apply}_{\text{ST}} &:: (\text{IsNat } ix, \text{EqHO } \kappa) \Rightarrow \text{Patch}_{\text{ST}} \ \kappa \text{ codes } ix \rightarrow \text{Fix } \kappa \text{ codes } ix \rightarrow \text{Maybe } (\text{Fix } \kappa \text{ codes } ix) \\
\text{apply}_{\text{ST}} &= \text{applyAl}\mu
\end{aligned}$$

An easily overlooked property of our patch definition is that the destination values it computes are guaranteed to be type-correct *by construction*. This is unlike the line-based or untyped approaches (which may generate ill-formed values) and similar to earlier results on type-safe differences [55].

FIGURE 4.2: A value of type $Patch_{ST}$ with its graphical representation.

4.2 MERGING PATCHES

The patches encoded in the $Patch_{ST}$ type clearly identify a prefix of constructors copied from the root of a tree up until the location of the changes and any insertion or deletions that might happen along the way. Moreover, since these patches also mirror the tree structure of the data in question, it becomes quite natural to identify separate changes. For example, if one change works on the left subtree of the root, and another on the right, they are clearly disjoint and can be merged. Finally, the explicit representation of insertions and deletions at the fixpoint level gives us a simple global alignment for our synchronizer.

In this section we discuss a simple merging algorithm, which reconciles changes from two different patches whenever these are *non-interfering*, for example, as in Figure 4.3. We call non-interfering patches *disjoint*, as they operate on separate parts of a tree.

A positive aspect of the $Patch_{ST}$ approach in comparison with a purely edit-scripts based approach is the significantly simpler merge function. This is due to $Patch_{ST}$ having clear homogeneous sections. Consequently, the type of the merge function is simple and reflects the fact that we expect a patch that operates over the values of the same type as a result:

$$merge :: Patch_{ST} \kappa \text{ codes } ix \rightarrow Patch_{ST} \kappa \text{ codes } ix \rightarrow Maybe (Patch_{ST} \kappa \text{ codes } ix)$$

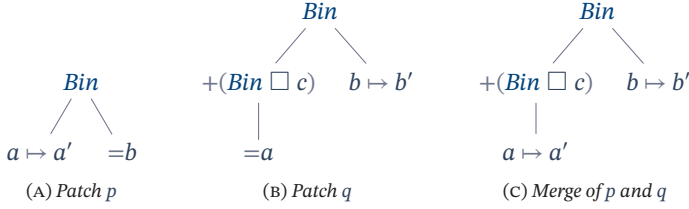


FIGURE 4.3: A simple example of mergeable patches.

A call to *merge*, in Haskell, returns *Nothing* if the patches have non-disjoint changes, that is, if both patches want to change the *same part* of the source tree.

Prior to prototyping *stdiff* in Haskell, we already had a working model of *stdiff* in Agda [74], which was created with the goal of proving that the merging algorithm would respect locality. In our Agda model, we have divided the merge function and the notion of disjointness, which yields a total merge function for the subset of disjoint patches:

$$\text{merge} : (p \ q : \text{Patch } \kappa \text{ codes } ix) \rightarrow \text{Disjoint } p \ q \rightarrow \text{Patch } \kappa \text{ codes } ix$$

A value of type *Disjoint* $p \ q$ corresponds to a proof that p and q change different parts of the source tree and is a symmetric relation – that is, *Disjoint* $p \ q$ iff *Disjoint* $q \ p$. This separation makes reasoning about the merge function much easier. In fact, we have proven that the merge function over regular datatypes commutes. A simplified statement of our theorem is given below:

$$\begin{aligned} \text{merge-commutes} & : (p \ q : \text{Patch } \kappa \text{ codes } ix) \\ & \rightarrow (\text{hyp} : \text{Disjoint } p \ q) \\ & \rightarrow \text{apply } (\text{merge } p \ q \ \text{hyp}) \circ q \equiv \text{apply } (\text{merge } q \ p \ (\text{sym } \text{hyp})) \circ p \end{aligned}$$

It is also worth noting that encoding the *merge* to be applied to the divergent replicas instead of the common ancestor – *residual-like* approach to merging (Section 2.1.4) – is instrumental to write a concise property and, consequently, prove the result. A merge function that applies to the common ancestor would probably require a much more convoluted encoding of *merge-commutes* above.

In a Haskell development, however, it is simpler to rely on the *Maybe* monad for disjointness. In fact, we define disjointness as whether or not merge returns a *Just*:

$$\begin{aligned} \text{disjoint} & :: \text{Patch } \kappa \text{ codes } ix \rightarrow \text{Patch } \kappa \text{ codes } ix \rightarrow \text{Bool} \\ \text{disjoint } p \ q & = \text{maybe } (\text{const } \text{True}) \ \text{False} \ (\text{merge } p \ q) \end{aligned}$$

The definition of the *merge* function is given in its entirety in Figure 4.4, but we discuss some interesting cases inline next. For example, when one change deletes a constructor but the other performs a change within said constructor we must check that they operate over *the same* constructor. When that is the case, we must go ahead and ensure the deletion context, *ctx*, and the changes in the product of atoms, *at*, are compatible.

$$\text{merge } (\text{Del } c_1 \text{ ctx}) (\text{Spn } (\text{SCns } c_2 \text{ at})) = \text{testEquality } c_1 \text{ } c_2 \gg \lambda \text{Refl} \rightarrow \text{mergeCtxAt } \text{ctx} \text{ } \text{at}$$

A (deletion) context is disjoint from a list of atoms if the patch in the hole of the context returns the same type of element as the patch on the product of patches and they are both disjoint. Moreover, the rest of the product of patches must consist in identity patches. Otherwise, we risk deleting newly introduced information.

```
mergeCtxAt :: DelCtx κ codes iy xs → NP (At κ codes) xs → Maybe (Alu κ codes ix iy)
mergeCtxAt (H (AlmuMin almu') rest) (AtFix almu × xs) = do
  Refl ← testEquality (almuDest almu) (almuDest almu')
  x ← mergeAlmu almu' almu
  guard (and $ elimNP identityAt xs)
  pure x
mergeCtxAt (T at ctx) (x × xs) = guard (identityAt x) >> mergeCtxAt ctx xs
```

The *testEquality* is there to ensure the patches to be merged are producing the same element of the mutually recursive family. This is one of the two places where we need these checks when adapting our Agda model to work over mutually recursive types. The second adaptation is shown shortly.

The *mergeAtCtx* function, dual to *mergeCtxAt*, merges a *NP (At κ codes) xs* and a *DelCtx κ codes iy xs* into a *Maybe (DelCtx κ codes iy xs)*, essentially preserving the *T at* it finds on the recursive calls. Another interesting case happens on one of the *mergeSpine* cases, whose full implementation can be seen in Figure 4.5. The *SChg* over *SCns* case must ensure we are working over the same element of the mutually recursive family, with a *testEquality ix iy*. This is the second place where we need to adapt the code in the Agda repository to work over mutually recursive types.

```
mergeSpine :: SNat ix → SNat iy
mergeSpine → Spine κ codes (Lkup ix codes) (Lkup iy codes)
mergeSpine → Spine κ codes (Lkup ix codes) (Lkup iy codes)
mergeSpine → Maybe (Spine κ codes (Lkup ix codes) (Lkup iy codes))
mergeSpine ix iy (SChg cx cy al) (SCns cz zs) = do
  Refl ← testEquality ix iy
  Refl ← testEquality cx cz
  SCns cy <$> mergeAlAt al zs
```

```

-- Non-disjoint recursive spines:
merge (Ins _ _) (Ins _ _) = Nothing
merge (Spn (SCHg _ _)) (Del _ _) = Nothing
merge (Del _ _) (Spn (Schg _ _)) = Nothing
merge (Del _ _) (Del _ _) = Nothing

-- Obviously disjoint recursive spines:
merge (Spn Scp) (Del c2 s2) = Just (Del c2 s2)
merge (Del c1 s2) (Spn Scp) = Just (Spn Scp)

-- Spines might be disjoint from spines and deletions:
merge (Spn s1) (Spn s2)
  = Spn <$> mergeSpine (getSNat (Proxy@ix)) (getSNat (Proxy@iy)) s1 s2
merge (Spn (SCns c1 at1)) (Del c2 s2)
  = Del c1 <$> mergeAtCtx at1 s2
merge (Del c1 s1) (Spn (SCns c2 at2))
  = do Refl ← testEquality c1 c2 -- disjoint if same constructor
    mergeCtxAt s1 at2

-- Insertions are disjoint from anything except insertions.
-- Overall disjointness does depend on the recursive parts, though.
merge (Ins c1 s1) (Spn s2) = Spn ∘ SCns c1 <$> mergeCtxAlmu s1 (Spn s2)
merge (Ins c1 s1) (Del c2 s2) = Spn ∘ SCns c1 <$> mergeCtxAlmu s1 (Del c2 s2)
merge (Spn s1) (Ins c2 s2) = Ins c2 <$> (mergeAlmuCtx (Spn s1) s2)
merge (Del c1 s1) (Ins c2 s2) = Ins c2 <$> (mergeAlmuCtx (Del c1 s1) s2)

```

FIGURE 4.4: Definition of merge.

```

-- Non-disjoint spines:
mergeSpine _ _ (SCHg _ _ _) (SCHg _ _ _) = Nothing

-- Obviously disjoint spines:
mergeSpine _ _ Scp s = Just s
mergeSpine _ _ s Scp = Just Scp

-- Disjointness depends on recursive parts:
mergeSpine _ _ (SCns cx xs) (SCns cy ys) = do Refl ← testEquality cx cy
    SCns cx <$> mergeAts xs ys
mergeSpine _ _ (SCns cx xs) (SCHg cy cz al) = do Refl ← testEquality cx cy
    SCHg cy cz <$> mergeAtAl xs al
mergeSpine ix iy (SCHg cx cy al) (SCns cz zs) = do Refl ← testEquality ix iy
    Refl ← testEquality cx cz
    SCns cy <$> mergeAlAt al zs

```

FIGURE 4.5: Definition of mergeSpine.

4.3 COMPUTING $Patch_{ST}$

In the previous sections, we have devised a typed representation for differences. We have seen that this representation is interesting in and of itself: being richly-structured and typed, it can be thought of as a non-trivial domain specific language whose denotation is given by the application function. Moreover, we have seen how to merge two disjoint differences. However, as programmers, we are mainly interested in *computing* patches from a source and a destination. Unfortunately, however, this is where the good news stops. Computing a value of type $Patch_{ST}$ is computationally expensive and represents one of the main downsides of this approach.

In this section we explore our attempts at computing differences with the `stdiff` framework. We start by outlining a nondeterministic specification of an algorithm for computing a $Patch_{ST}$, in Section 4.3.1. We then provide example algorithms that implemented said specification in Section 4.3.2. All these approaches, however, need to make choices. Moreover, the rich structure of $Patch_{ST}$ makes a memoized algorithm much more difficult to write. Consequently, computing a $Patch_{ST}$ will always be a computationally inefficient process, rendering it unusable in practice.

4.3.1 NAIVE ENUMERATION

The simplest option for computing a patch that transforms a tree x into y is enumerating all possible patches and filtering out those with the smallest *cost*, for some *cost* metric. In this section, we will write a naive enumeration engine for $Patch_{ST}$ and argue that regardless of the *cost* notion, the state space explodes quickly and becomes intractable.

The enumeration follows the Agda model [74] closely and is not very surprising. Nevertheless, it does act as a good specification for a better implementation later. Just like for the linear case, the changes that can transform two values x and y of a given mutually recursive family into one another are the deletion of a constructor from x , the insertion of a constructor from y or changing the constructor of x into the one from y , as witnessed by the `enumAlμ` function below.

```
enumAlμ :: Fix ki codes ix → Fix ki codes iy → [Alμ ki codes ix iy]
enumAlμ x y = enumDel (sop $ unFix x) y
              <|> enumIns x (sop $ unFix y)
              <|> Spn <$> enumSpn (snatFixIdx x) (snatFixIdx y)
                  (unFix x) (unFix y)
```

where

```
enumDel (Tag c p) y₀ = Del c <$> enumDelCtx p y₀
enumIns x₀ (Tag c p) = Ins c <$> enumInsCtx x₀ p
```

Enumerating all the patches from a deletion context of a given product p against some fixpoint y consists of enumerating the patches that transform all of the fields of p into y . The handling of insertion contexts is analogous, hence it is omitted here. Recall that the *AlmuMin*, below, is used to flag the resulting context as a deletion context.

```
enumDelCtx :: PoA ki (Fix ki codes) prod → Fix ki codes iy → [DelCtx ki codes iy prod]
enumDelCtx Nil _ = []
enumDelCtx (NAK x × xs) f = T (NAK x) <$> enumDelCtx xs f
enumDelCtx (NAI x × xs) f = (flip H xs ∘ AlmuMin) <$> enumAlμ x f
                           <|> T (NAI x) <$> enumDelCtx xs f
```

Next we look into enumerating the spines between x and y , that is, changes to the coproduct structure from x to y . Unlike our Agda model, we need to know over which element of the mutually recursive family we are operating. This will dictate which constructors from *Spine* we are allowed to use. We gather this information through two auxiliary *SNat* parameters. The choice of which spine constructor to use is deterministic, that is, each case is uniquely determined by a *Spine* constructor.

```
enumSpn :: SNat ix → SNat iy
        → Rep ki (Fix ki codes) (Lkup ix codes)
        → Rep ki (Fix ki codes) (Lkup iy codes)
        → [Spine ki codes (Lkup ix codes) (Lkup iy codes)]
enumSpn six siy x y =
  let Tag cx px = sop x
      Tag cy py = sop y
  in case testEquality six siy of
    Nothing → SChg cx cy <$> enumAl px py
    Just Refl → case testEquality cx cy of
      Nothing → SChg cx cy <$> enumAl px py
      Just Refl → if eqHO px py
        then return Scp
        else SCns cx <$> mapNPM (uncurry' enumAt) (zipNP px py)
```

Enumerating atoms, *enumAt*, is trivial. Atoms are either opaque types or recursive positions. Opaque types are handled by *TrivialK* and recursive positions are handled recursively by *enumAlμ*.

```
enumAt :: NA ki (Fix ki codes) at → NA ki (Fix ki codes) at → [At ki codes at]
enumAt (NAI x) (NAI y) = AtFix <$> enumAlμ x y
enumAt (NAK x) (NAK y) = return $ AtSet (Trivial x y)
```

Finally, product alignment is analogous to the longest common subsequence, except that we must make sure that we only synchronize atoms with *AX* if they have the same

type. The *enumAl* below illustrates the non-deterministic enumeration of alignments over two products-of-atoms.

```

enumAl :: PoA ki (Fix ki codes) p1 → PoA ki (Fix ki codes) p2 → [Al ki codes p1 p2]
enumAl Nil Nil = return A0
enumAl (x × xs) Nil = ADel x <$> enumAl xs Nil
enumAl Nil (y × ys) = AIns y <$> enumAl Nil ys
enumAl (x × xs) (y × ys) = (ADel x <$> enumAl xs (y × ys))
                           <|> (AIns y <$> enumAl (x × xs) ys)
                           <|> case testEquality x y of
                               Just Refl → AX <$> (enumAt x y) <*> enumAl xs ys
                               Nothing → mzero

```

From the definitions of *enumAlμ* and *enumAl*, it is clear why this algorithm explodes and becomes intractable. In *enumAlμ* we must choose between inserting, deleting or copying a recursive constructor. In case we chose to copy a constructor, we then might call *enumAl*, where we must choose between inserting, deleting or copying fields of constructors. We must enumerate these options for virtually each pair of constructors in the source and destination trees.

4.3.2 TRANSLATING FROM GDIFF

Since enumerating all possible patches and then filtering a chosen one is time consuming and requires a complex notion of cost over $Patch_{ST}$, it was clear we should be pursuing better algorithms for our *diff* function. We have attempted two similar approaches to filter the uninteresting patches out and optimize the search space.

A first idea, which arose in collaboration with Pierre-Evariste Dagand (private communication), was to use the already existing UNIX *diff* tool as some sort of *oracle*. That is, we should only consider inserting and deleting elements that fall on lines marked as such by UNIX *diff*. This idea was translated into Haskell by Garuffi [36], but the performance was still very poor and computing the $Patch_{ST}$ of two real-world Clojure files still required several minutes.

From Garuffi’s experiments [36] we learnt that simply restricting the search space was not sufficient. Besides the complexity introduced by arbitrary heuristics, using the UNIX *diff* to flag elements of the AST was still too coarse. For one, the UNIX *diff* can insert and delete the same line in some situations. Secondly, many elements of the AST may fall on the same line.

The second option is related, but instead of using a line-based oracle, we can use *gdiff* (Section 3.1.4) as the oracle, enabling us to annotate every node of the source and destination trees with information about whether that node was copied or not. This strategy was translated into Haskell by Van Putten [92] as part of his MSc work. The gist

of it is that we can use annotated fixpoints to tag each constructor of a tree with added information. In this case, we are interested in whether this node would be copied or not by `gdifff`:

data *Ann* = *Modify* | *Copy*

A *Modify* annotation corresponds to a deletion or insertion depending on whether it is the source or destination tree respectively. Recall that an edit-script produced by `gdifff` has type *ES* κ codes *xs ys*, where *xs* is the list of types of the source trees and *ys* is the list of types of the destination trees. The definition of *ES* – introduced in Section 3.1.4 – is repeated below.

data *ES* κ codes :: [*Atom* *kon*] → [*Atom* *kon*] → * **where**
ES0 :: *ES* κ codes '[]' '[]'
Ins :: *Cof* κ codes *a t* → *ES* κ codes *i* (*t* : $\#$: *j*) → *ES* κ codes *i* (*a* ' \prime : *j*)
Del :: *Cof* κ codes *a t* → *ES* κ codes (*t* : $\#$: *i*) *j* → *ES* κ codes (*a* ' \prime : *i*) *j*
Cpy :: *Cof* κ codes *a t* → *ES* κ codes (*t* : $\#$: *i*) (*t* : $\#$: *j*) → *ES* κ codes (*a* ' \prime : *i*) (*a* ' \prime : *j*)

Given a value of type *ES* κ codes *xs ys*, we have information about which constructors of the trees in *NP* (*NA* κ (*Fix* κ codes)) *xs* should be copied. Our objective then is to annotate the trees with this very information. This is done by the *annSrc* and *annDst* functions. We will only look at *annSrc*, the definition of *annDst* is symmetric.

Annotating the source forest with a given edit-script consists in matching which constructors present in the forest correspond to a copy and which correspond to a deletion. The insertions in the edit-script concern the destination forest only. The *annSrc* function, below, does exactly that, proceeding by induction on the edit-script.

annSrc :: *NP* (*NA* κ (*Fix* κ codes)) *xs* → *ES* κ codes *xs ys*
 → *NP* (*NA* κ (*FixAnn* κ codes (*Const Ann*))) *xs*
annSrc xs ES0 = *Nil*
annSrc Nil = *Nil*
annSrc xs (Ins c es) = *annSrc' xs es*
annSrc (x \times xs) (Del c es) = **let** *poa* = *fromJust* \$ *matchCof* *c x*
 in *insCofAnn c (Const Modify) (annSrc' (appendNP poa xs) es)*
annSrc' (x \times xs) (Cpy _ c es) = **let** *poa* = *fromJust* \$ *matchCof* *c x*
 in *insCofAnn c (Const Copy) (annSrc' (appendNP poa xs) es)*

The deterministic diff function for *Al μ* starts by checking the annotations present at the root of its argument trees. In case both are copies, we start with a spine. If at least one of them is not a copy we insert or delete the constructor not flagged as a copy. We must guard for the case that there exists a copy in the inserted or deleted subtree. In case that does not hold, we would not be able to choose an argument of the inserted or deleted constructor to continue diffing against, in *diffCtx*. When there are no more

copies to be performed, we just return a *stiff* patch, which deletes the entire source and inserts the entire destination tree.

```

diffAlmu :: FixAnn  $\kappa$  codes (Const Ann) ix  $\rightarrow$  FixAnn  $\kappa$  codes (Const Ann) iy
 $\rightarrow$  Almu  $\kappa$  codes ix iy
diffAlmu x@(FixAnn ann1 rep1) y@(FixAnn ann2 rep2) =
  case (getAnn ann1, getAnn ann2) of
    (Copy, Copy)  $\rightarrow$  Spn (diffSpine (getSNat $ Proxy@ix)
                                   (getSNat $ Proxy@iy)
                                   rep1 rep2)
    (Copy, Modify)  $\rightarrow$  if hasCopies y then diffIns x rep2
                        else stiffAlmu (forgetAnn x) (forgetAnn y)
    (Modify, Copy)  $\rightarrow$  if hasCopies x then diffDel rep1 y
                        else stiffAlmu (forgetAnn x) (forgetAnn y)
    (Modify, Modify)  $\rightarrow$  if hasCopies x then diffDel rep1 y
                        else stiffAlmu (forgetAnn x) (forgetAnn y)
  where
    diffIns x rep = case sop rep of Tag c ys  $\rightarrow$  Ins c (diffCtx CtxIns x ys)
    diffDel rep y = case sop rep of Tag c xs  $\rightarrow$  Del c (diffCtx CtxDel y xs)

```

The *diffCtx* function selects an element of a product to continue diffing against. We naturally select the element that has the most constructors marked for copy as the element we continue diffing against. The other fields of the product are placed on the *rigid* part of the context, that is, the trees that will be deleted or inserted in their entirety, without sharing any of their subtrees.

```

diffCtx :: InsOrDel  $\kappa$  codes p  $\rightarrow$  FixAnn  $\kappa$  codes (Const Ann) ix
 $\rightarrow$  NP (NA  $\kappa$  (FixAnn  $\kappa$  codes (Const Ann))) xs
 $\rightarrow$  Ctx  $\kappa$  codes p ix xs

```

The other functions for translating two *FixAnn κ codes (Const Ann) ix* into a *Patch_{ST}* are straightforward and follow a similar reasoning process: extract the annotations and defer copies until both source and destination annotation flag a copy.

This version of the *diff* function runs in $\mathcal{O}(n^2)$ time, where n is the the number of constructors in the bigger input tree. Although orders of magnitude better than naive enumeration or using the UNIX *diff* as an oracle, a quadratic algorithm is still not practical, particularly when n tends to be large – real-world source files have tens of thousands abstract syntax elements.

4.4 DISCUSSION

With `stdiff` we learned that the difficulties of edit-script based approaches are not due, exclusively, to using linear data to represent transformations to tree structured data. Another important aspect that we unknowingly overlooked, and ultimately did lead to a prohibitively expensive *diff* function, was the necessity to choose a single copy opportunity. This happens whenever a subtree could be copied in two or more different ways, and, in tree differencing this occurs often. For example, think of all the places that a call to a logging function *log err msg* could be copied in a large source-file; or all of the *+1* expressions.

The *Patch_{ST}* datatype has many interesting aspects that deserve some mention. First, by being globally synchronized – that is, explicit insertions and deletions with one hole – these patches are easy to merge. Moreover, we have seen that it is possible, and desirable, to encode patches as homogeneous types: a patch transforms two values of the same member of the mutually recursive family.

In conclusion, lacking an efficient *diff* algorithm meant that `stdiff` was an important step leading to new insights, but unfortunately was not worth pursuing further. This meant that a number of interesting topics such as the algebra of *Patch_{ST}* and the notion of cost for *Patch_{ST}* were abandoned indefinitely.



PATTERN-EXPRESSION PATCHES

The `stdiff` approach gave us a first representation of tree-structured patches over tree-structured data but was still deeply connected to edit-scripts: subtrees could only be copied once and could not be permuted. This means we still suffered from ambiguous patches, and, consequently, a computationally expensive *diff* algorithm. Overcoming the drawback of ambiguity requires a shift in perspective and abandoning edit-script based differencing algorithms. In this section we will explore the `hdiff` approach, where patches allow for trees to be arbitrarily permuted, duplicated or contracted (contractions are dual to duplications).

Classical tree differencing algorithms start with the computation of tree matchings (Section 2.1.2), which identify the subtrees that should be copied. These tree matchings, however, must be restricted to order-preserving partial injections to be efficiently translated to edit-scripts later. The `hdiff` approach never translates to edit-scripts, which means the tree matchings we compute are not subject to *any* restrictions. In fact, `hdiff` uses these unrestricted tree matchings as *the patch*, instead of translating them *into* a patch.

Suppose we want to describe a change that modifies the left element of a binary tree. If we had the full Haskell programming language available as the patch language, we could write something similar to function `c`, in Figure 5.1(A). Observing the function `c` we see it has a clear domain – a set of *Trees* that when applied to `c` yields a *Just* – which is specified by the pattern and guards. Then, for each tree in the domain we compute a corresponding tree in the codomain. The new tree is obtained from the old tree by replacing the `10` by `42` in-place. Closely inspecting this definition, we can interpret the matching of the pattern as a *deletion* phase and the construction of the resulting tree as a *insertion* phase. The `hdiff` approach represents the change in `c` exactly as that: a pattern

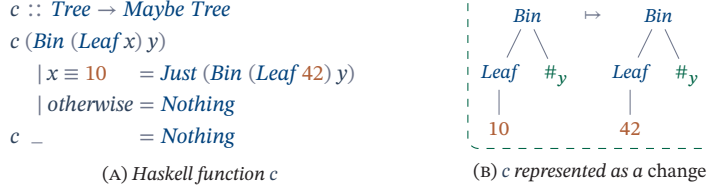


FIGURE 5.1: Haskell function and its graphical representation as a change. The change here modifies the left child of a binary node. Notation $\#_y$ is used to indicate y is a metavariable.

and a expression. Essentially, we write c as $\text{Chg } (\text{Bin } (\text{Leaf } 10) y) (\text{Bin } (\text{Leaf } 42) y)$ – represented graphically as in Figure 5.1(B).

With the added expressivity of referring to subtrees with metavariables we can represent more transformations than before. Take, for example, the change that swaps two subtrees – which cannot be written using an edit-script based approach – given by $\text{Chg } (\text{Bin } x y) (\text{Bin } y x)$. Another helpful consequence of our design is that we effectively bypass the *choice* phase of the algorithm. When computing the differences between Bin Leaf Leaf and Leaf , for example, we do not have to chose one Leaf to copy because we can copy both with the help of a contraction operation, with a change such as: $\text{Chg } (\text{Bin } x x) x$. This aspect is crucial and enables us to write a linear *diff* algorithm.

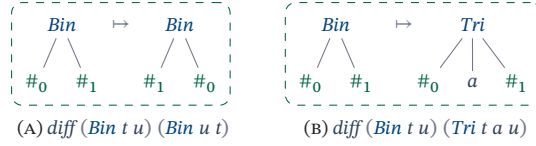
In this chapter we explore the representation and computation aspects of `hdiff`. The big shift in paradigm of `hdiff` also requires a more careful look into the metatheory and nuances of the algorithm, which were not present in our original paper [76]. The material in this chapter is developed from our ICFP’19 publication [76], shifting to the `generics-simplistic` library.

5.1 CHANGES

5.1.1 A CONCRETE EXAMPLE

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

```
data Tree = Leaf Int | Bin Tree Tree | Tri Tree Tree Tree
```

FIGURE 5.2: Illustration of two changes. Metavariables are denoted with $\#_x$.

The central concept of `hdiff` is the encoding of a *change*. Unlike previous work [55, 74, 50] which is based on tree-edit-distance [16] and hence uses only insertions, deletions and copies of the constructors encountered during the preorder traversal of a tree (Section 3.1.4), we go a step further. We explicitly model permutations, duplications and contractions of subtrees within our notion of *change*, where contraction here denotes the partial inverse of a duplication. The representation of a *change* between two values of type *Tree*, 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, *TreeC* φ , enables us to annotate a value of *Tree* with holes of type φ . Therefore, *TreeC Metavar* represents the type of *Tree* 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 *TreeC* as a *context*. For now, the metavariables will be simple *Int* values but later on they will need to carry additional information.

```

type Metavar = Int
data TreeC  $\varphi$  = Hole  $\varphi$ 
              | LeafC Int
              | BinC TreeC TreeC
              | TriC TreeC TreeC TreeC

```

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 Change  $\varphi$  = (TreeC  $\varphi$ , TreeC  $\varphi$ )

```

The change that transforms *Bin t u* into *Bin u t*, for example, is represented by a pair of *TreeC*, (*BinC* (*Hole* 0) (*Hole* 1), *BinC* (*Hole* 1) (*Hole* 0)), as seen in Figure 5.2. This change works on any tree built using the *Bin* 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.

5.1.1.1 APPLYING CHANGES

Applying a change to a tree is done by unifying the metavariables in the deletion context with said tree, and later instantiating the insertion context with the obtained substitution. Later on, when we come to the generic setting, we will write the application function using syntactic unification [94]. For this concrete example, we will continue with the definition below.

```
chgApply :: Change Metavar → Tree → Maybe Tree
chgApply (d, i) x = del d x ≧ 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, i.e., the contexts are non-linear. 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. Given a deletion context ctx and source tree t , the *del* function tries to associate all the metavariables in the context with a subtree of the input *tree*. This can be done with standard unification algorithms, as will be the case in the generic setting. Here, however, we use a simple auxiliary function to do so.

```
del :: TreeC Metavar → Tree → Maybe (Map Metavar Tree)
del ctx t = go ctx t empty
```

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

```
go :: TreeC → Tree → Map Metavar Tree → Maybe (Map Metavar Tree)
go (LeafC n) (Leaf n') m = guard (n ≡ n') ≧ return m
go (BinC x y) (Bin a b) m = go x a m ≧ go y b
go (TriC x y z) (Tri 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
```

Once we have computed the substitution that unifies ctx and t , above, we instantiate the variables in the insertion context with their respective values to obtain the resulting tree. The ins function, defined below, performs this instantiation and fails only if the change contains unbound variables.

```

ins :: TreeC Metavar → Map Metavar Tree → Maybe Tree
ins (LeafC n) m = return (Leaf n)
ins (BinC x y) m = Bin <$> ins x m <*> ins y m
ins (TriC x y z) m = Tri <$> ins x m <*> ins y m <*> ins z m
ins (Hole i) m = lookup i m

```

5.1.1.2 COMPUTING CHANGES

Next we will define the $chgTree$ function, which produces a change from a source and a destination. Intuitively, the $chgTree$ function should 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 to a function $wcs :: Tree \rightarrow Tree \rightarrow Tree \rightarrow Maybe Metavar$, 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 $wcs\ s\ d\ x \equiv wcs\ s\ d\ y \equiv Just\ j$, then $x \equiv y$. 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 5.1.1.3, we provide an efficient implementation. For now, assuming the oracle exists allows for a clear separation of concerns. The $chgTree$ function merely has to compute the deletion and insertion contexts, using said oracle – the inner workings of the oracle are abstracted away cleanly.

```

chgTree :: Tree → Tree → Change Metavar
chgTree s d = let f = wcs s d
              in (extract f s, extract f d)

```

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.

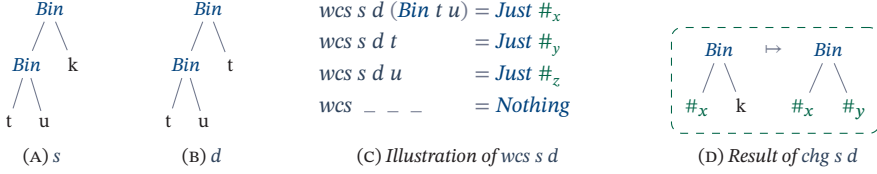


FIGURE 5.3: Context extraction must care to produce well-formed changes. The nested occurrence of t within $Bin\ t\ u$ here yields a change with an undefined variable on its insertion context.

If the oracle returns *Nothing*, we move the topmost constructor to the context being computed and recurse over the remaining subtrees.

```

extract :: (Tree → Maybe Metavar) → Tree → TreeC Metavar
extract o t = maybe (peel t) Hole (o t)
  where peel (Leaf n)   = LeafC n
        peel (Bin a b)  = BinC (extract o a) (extract o b)
        peel (Tri a b c) = TriC (extract o a) (extract o b) (extract o c)

```

Note that by adopting 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 *chgTree* function has a subtle bug, however. It does *not* produce correct changes, that is, it is not the case that $apply\ (chg\ s\ d)\ s \equiv Just\ d$ for all s and d . The problem can be observed when we pass a source and a destination tree where a common subtree occurs by itself but also as a subtree of another common subtree. Such situation is illustrated in Figure 5.3. In particular, the patch shown in Figure 5.3(D) cannot be applied since the deletion context does not instantiate the metavariable $\#_y$, which is required by the insertion context.

There are many ways to address the issue illustrated in Figure 5.3. We could replace $\#_y$ by t and ignore the sharing or we could replace $\#_x$ by $Bin\ \#_y\ \#_z$, pushing the metavariables to the leaves maximizing sharing. These would give rise to the changes shown in Figure 5.4. There is a clear dichotomy between wanting to maximize the spine but at the same time wanting to copy the larger trees, closer to the root. On the one hand, copies closer to the root are intuitively easier to merge and less sharing means it is easier to isolate changes to separate parts of the tree. On the other hand, sharing as much as possible might capture the change being represented more closely.

A third, perhaps less intuitive, solution to the problem in Figure 5.3 is to only shares uniquely occurring subtrees, effectively simulating the UNIX *diff* with the *patience*

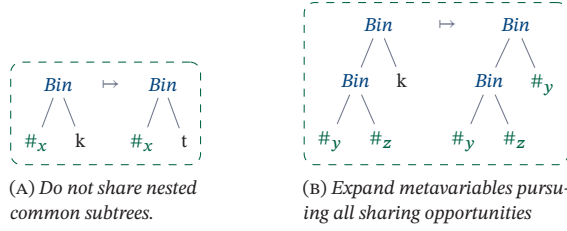


FIGURE 5.4: Two potential solutions to the problem of nested common subtrees, illustrated in Figure 5.3.

option, which only copies uniquely occurring lines. In fact, to make this easy to experiment with, we will parameterize our final *extract* with which *context extraction mode* should be used to computing changes.

```
data ExtractionMode = NoNested
                    | ProperShare
                    | Patience
```

The *NoNested* mode will forget sharing in favor of copying larger subtrees. It would drop the sharing of *t* producing Figure 5.4(A). The *ProperShare* mode is the opposite. It would produce Figure 5.4(B). Finally, *Patience* only share subtrees that occur only once in the source and once in the destination. For the inputs in Figure 5.3, extracting contexts under *Patience* mode would produce the same result as *NoNested*, but they are not the same in general. In fact, Figure 5.5 illustrates the changes that would be extracted following each *ExtractionMode* for the same source and destination.

In short, the *extract* function receives the *sharing map* and extracts the deletion and insertion context making up the change, caring that the produced change is well-scoped. We will give the final *extract* function when we get to its generic implementation. For the time being, let us move on to the intuition behind computing the *wcs* function efficiently for the concrete case of the *Tree* datatype.

5.1.1.3 DEFINING THE ORACLE FOR *Tree*

In order to have a working version of our differencing algorithm for *Tree* we must provide the *wcs* implementation, with type $Tree \rightarrow Tree \rightarrow Tree \rightarrow \text{Maybe Metavar}$. Given a fixed *s* and *d*, *wcs s d x* returns *Just i* if *x* is the *i*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

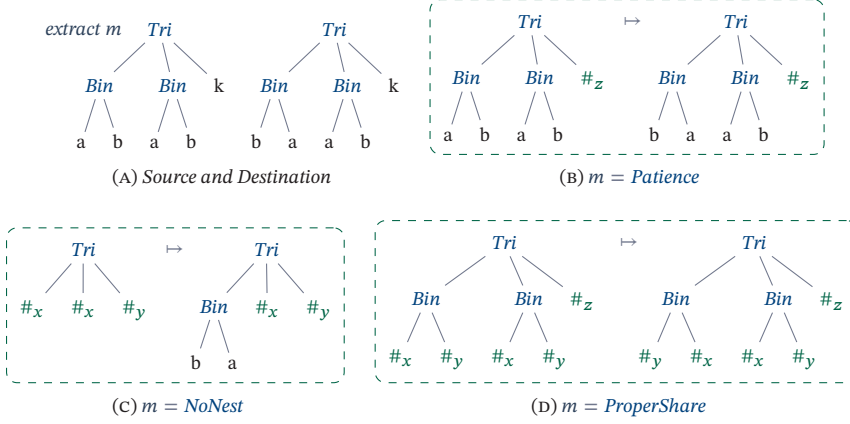


FIGURE 5.5: Different extraction methods for the same pair or trees.

the subtrees in s and d , and then searches for the subtree x the resulting list. Enumerating all the subtrees of any *Tree* is straightforward:

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

It is now easy 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.

$\text{wcs} :: \text{Tree} \rightarrow \text{Tree} \rightarrow \text{Tree} \rightarrow \text{Maybe Metavar}$
 $\text{wcs } s \ d \ x = \text{elemIndex } x \ (\text{subtrees } s \cap \text{subtrees } d)$

This specification, however, is not particularly efficient. The inefficiency comes from two places: computing $\text{subtrees } s \cap \text{subtrees } d$ is expensive but could be cached, but the call to $\text{elemIndex } x$ will be repeated for each subtree of s and d when extracting the contexts. This means that the overall algorithm will be close to exponential. In fact, given how often we need to call wcs , each call to must run in amortized constant time if we want our algorithm to be efficient.

Defining $\text{wcs } s \ d$ efficiently consists, firstly, of computing a set of trees which contains the subtrees of s and d , and secondly, in being able to efficiently query this set for membership. Symbolic manipulation software, such as Computer Algebra Systems, perform similar computations frequently and their performance is just as important. These systems often rely on a technique known as *hash-consing* [39, 33], which is part of the canon of programming folklore. Hash-consing arises as a means of *maximal sharing*

of subtrees in memory and constant time comparison – two trees are equal if they are stored in the same memory location – but it is by far not limited to it. We will be using a variant of *hash-consing* to define *wcs s d*.

To efficiently compare trees for equality we will be using cryptographic hash functions [69] 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* [70]. 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 Tree* **where**

t \equiv *u* = *merkleRoot t* \equiv *merkleRoot u*

The definition of *merkleRoot* function is straightforward. It is important that we use the *merkleRoot* of the parts of a *Tree* 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 :: Tree → Digest
merkleRoot (LeafH n) = hash (concat [ "1", encode n ])
merkleRoot (Bin x y) = hash (concat [ "2", merkleRoot x, merkleRoot y ])
merkleRoot (Tri x y z) = hash (concat [ "3", 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 [69]. Nonetheless, it would be 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 comparison, but the (\equiv) definition comparing merkle roots is still linear as it must recompute the *merkleRoot* on every comparison. We fix this by caching the hash associated with every node of a *Tree*. This is done by the *decorate* function, illustrated Figure 5.6.

```
type TreeH = (TreeH', Digest)
data TreeH' = LeafH Int
              | BinH TreeH TreeH
              | TriH TreeH TreeH TreeH
decorate :: Tree → TreeH
```

We omit the implementation of *decorate* for now, even if it is straightforward. Moreover, a generic version is introduced in Section 5.1.4. The *TreeH* datatype carries round the *merkleRoot* of its first component, hence, enabling us to define (\equiv) in constant time.

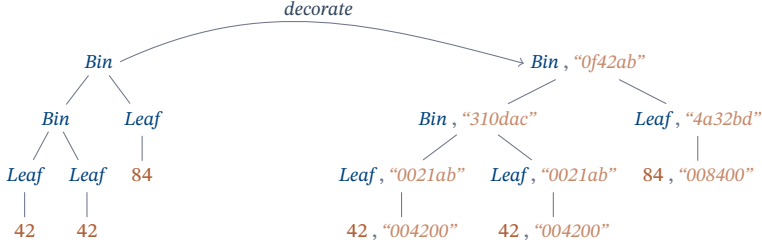


FIGURE 5.6: Example of decorating a tree with hashes, through the *decorate* function.

The second source of inefficiency, enumerating all possible subtrees, can be addressed by choosing a better data structure. To check whether a tree x is a subtree of a fixed s and d , it suffices 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 $[\text{Word}]$, the optimal choice for such “database” is a Trie [19] mapping a $[\text{Word}]$ 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 *Tree* type. The *wcs* oracle will now receive *TreeH*, i.e., trees annotated with their merkle roots at every node, and will populate the “database” of common digests.

```
wcs :: TreeH → TreeH → TreeH → Maybe Metavar
wcs s d = lookup (mkTrie s ∩ mkTrie d) ∘ merkleRoot
  where (∩)    :: Trie k v → Trie k u → Trie k v -- key inter., keep left values.
         lookup :: Trie k v → [k]      → Maybe v
         mkTrie :: TreeH  → 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 note that the efficiency of the algorithm comes from the novel representation of patches combined with an amortized constant time *wcs* function. Without being able to duplicate or permute subtrees, the algorithm would have to backtrack in a number of situations.

5.1.2 REPRESENTING CHANGES GENERICALLY

Having seen how *TreeC* played a crucial role in defining changes for the *Tree* datatype, we continue with its generic implementation. In this section, we generalize the construction of *contexts* to any datatype supported by the *generics-simplicistic* library.

Recall that a *context* over a datatype T is just a value of T augmented with an additional constructor used to represent *holes*. This can be done with the *free monad* construction provided by the `generics-simplistic` library: `HolesAnn κ fam ann h` datatype (Section 3.2.3) is a free monad in h . We recall its definition ignoring annotations below.

```
data Holes  $\kappa$  fam h a where
  Hole  :: h a  $\rightarrow$  Holes  $\kappa$  fam h a
  Prim  :: (PrimCnstr  $\kappa$  fam a)  $\Rightarrow$  a  $\rightarrow$  Holes  $\kappa$  fam h a
  Roll  :: (CompoundCnstr  $\kappa$  fam a)  $\Rightarrow$  SRep (Holes  $\kappa$  fam h) (Rep a)  $\rightarrow$  Holes  $\kappa$  fam h a
```

The *TreeC Metavar* datatype, defined in Section 5.1.1 to represent a value of type *Tree* augmented with metavariables is isomorphic to `Holes '[Int] '[Tree] (Const Int)`. Abstracting over the specific family for *Tree*, the datatype `Holes κ fam (Const Int)` gives a functor mapping an element of the family into its representation augmented with integers, which represent metavariables. But in this generic setting, it does not yet enable us to infer whether a metavariable matches over an opaque type or a recursive position, which will come to be important soon. Consequently, we will keep the information about whether the metavariable matches over an opaque value or not:

```
data Metavar  $\kappa$  fam at where
  # $\kappa$  :: (PrimCnstr  $\kappa$  fam at)
       $\Rightarrow$  Int  $\rightarrow$  Metavar  $\kappa$  fam at
  #fam :: (CompoundCnstr  $\kappa$  fam at)
       $\Rightarrow$  Int  $\rightarrow$  Metavar  $\kappa$  fam at
```

With *Metavar* above, we can always retrieve the *Int* identifying the metavar, with the `metavarGet` function, but we maintain all the type-level information we may need to inspect at run-time. The *HolesMV* datatype below is convenient since most of the times our *Holes* structures will contain metavariables.

```
metavarGet :: Metavar  $\kappa$  fam at  $\rightarrow$  Int
type HolesMV  $\kappa$  fam = Holes  $\kappa$  fam (Metavar  $\kappa$  fam)
```

A *change* consists of a pair of a deletion context and an insertion context for the same type. These contexts are values of the mutually recursive family in question, augmented with metavariables.

```
data Chg  $\kappa$  fam at = Chg {  $\cdot_{\text{del}}$  :: HolesMV  $\kappa$  fam at
                        ,  $\cdot_{\text{ins}}$  :: HolesMV  $\kappa$  fam at
                        }
```

Applying a generic change c to an element x consists in unifying x with c_{del} , yielding a substitution σ which can be applied to c_{ins} . This provides the usual denotational semantics of changes as partial functions.

```

chgApply :: (All Eq  $\kappa$ )  $\Rightarrow$   $\text{Chg } \kappa \text{ fam at} \rightarrow \text{SFix } \kappa \text{ fam at} \rightarrow \text{Maybe } (\text{SFix } \kappa \text{ fam at})$ 
chgApply ( $\text{Chg } d \ i$ )  $x$  = either (const Nothing) (holesMapM uninstantHole  $\circ$  flip substApply  $i$ )
                               (unify  $d$  (sfixToHoles  $x$ ))
where uninstantHole _ = error "uninstantiated hole: ( $\text{Chg } d \ i$ ) not well-scoped!"

```

In a call to `chgApply c x` , since x has no holes, a successful unification means σ assigns a term (no holes) for each metavariable in c_{del} . In turn, when applying σ to c_{ins} we must guarantee that every metavariable in c_{ins} gets substituted, yielding a term with no holes as a result. Attempting to apply a non-well-scoped change is a violation of the contract of `applyChg`. We throw an error in that case and distinguish it from a change c not being able to be applied to x because x is not an element of the domain of c . The `uninstantHole` above will be called in the precise situation where holes were left uninstantiated in `substApply σ c_{ins}` .

In general, we expect a value of type `Chg` to be well-scoped, that is, all the variables that are present in the insertion context must also occur on the deletion context, in Haskell:

```

vars      :: HolesMV  $\kappa$  fam at  $\rightarrow$  Map Int Arity
wellscored ::  $\text{Chg } \kappa \text{ fam at} \rightarrow \text{Bool}$ 
wellscored ( $\text{Chg } d \ i$ ) = keys (vars  $i$ )  $\equiv$  keys (vars  $d$ )

```

This definition of well-scoped changes was chosen due to its simplicity. If we know a given variable has arity 2, for example, then it *must* be a copy. Had we defined well-scoped changes by `keys (vars i) \subseteq keys (vars d)`, a variable with arity 2 could have both its occurrences in the deletion context. This would make the merging algorithm even more involved.

A `Chg` is very similar to a *tree matching* (Section 2.1.2) with less restrictions. In other words, it arbitrarily maps subtrees from the source to the destination. From an algebraic point of view, this already gives us a desirable structure, as we will explore next in Section 5.1.3. In fact, we argue that there is no need to translate the tree matching into an edit-script, like most traditional algorithms do. The tree matching should be used as the representation of change.

5.1.3 META THEORY

In this section we will look into how `Chg` admits a simple composition operation which makes a partial monoid. Through the remainder of this section we will assume changes

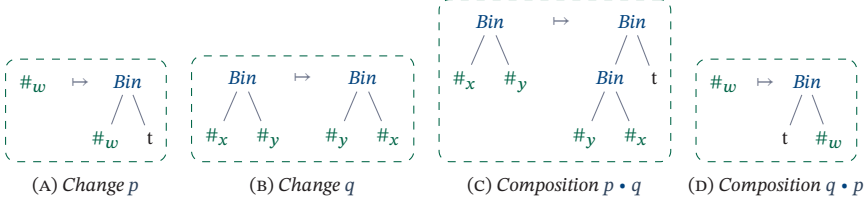


FIGURE 5.7: Example of change composition. The composition usually can be applied to less elements than its parts and is clearly not commutative.

have all been α -converted to never capture names and denote the application function of a change, $\text{applyChg } c$, as $\llbracket c \rrbracket$. We will also abuse notation and denote $\text{substApply } \sigma p$ by σp , whenever the context makes it clear that σ is a substitution. Finally, we will abide by the Barendregt convention [13] in our proofs and metatheory – that is, all changes that appear in some mathematical context have their bound variable names independent of each other, to put it differently, no two changes will accidentally share a variable name.

The composition of two changes, say, p after q , returns a change that maps a subset of the domain of q into a subset of the image of p . Figure 5.7, for example, illustrates two changes and their two different compositions. In the case of Figure 5.7 both $p \bullet q$ and $q \bullet p$ exist, but this is not the case generally. The composition of two changes $p \bullet q$ is defined if and only if the image of $\llbracket q \rrbracket$ has some elements in common with the domain of $\llbracket p \rrbracket$. In other words, when q_{ins} is unifiable with p_{del} . In fact, let $\sigma = \text{unify } q_{\text{ins}} p_{\text{del}}$, the composition $p \bullet q$ is given by $\text{Chg } (\sigma q_{\text{del}}) (\sigma p_{\text{ins}})$.

$(\bullet) :: \text{Chg } \kappa \text{ fam at} \rightarrow \text{Chg } \kappa \text{ fam at} \rightarrow \text{Maybe } (\text{Chg } \kappa \text{ fam at})$

$p \bullet q = \text{case unify } p_{\text{del}} q_{\text{ins}} \text{ of}$

Left $_ \rightarrow \text{Nothing}$

Right $\sigma \rightarrow \text{Just } (\text{Chg } (\sigma q_{\text{del}}) (\sigma p_{\text{ins}}))$

Note that it is inherent that purely structural composition of two changes p after q yields a change, $p \bullet q$, that potentially misses sharing opportunities. Imagine that p inserts a subtree t that was deleted by q . Our composition algorithm possesses no information that this t is to be treated as a copy. This also occurs in the edit-script universe: composing patches yields worse patches than recomputing differences. We can imagine that a more complicated composition algorithm might be able to recover the copies in those situations.

We do not particularly care whether composition produces *the best* change possible or not. We do not even have a notion of *best* at the moment. It is vital, however, that it produces a correct change. That is, the composition of two patches is indistinguishable from the composition of their application functions.

Lemma 5.1.1 (Composition Correct). *For any changes p and q and trees x and y aptly typed; we have $\llbracket p \bullet q \rrbracket x \equiv \text{Just } y$ if and only if $\exists z. \llbracket q \rrbracket x \equiv \text{Just } z \wedge \llbracket p \rrbracket z \equiv \text{Just } y$.*

Proof. **if.** Assuming $\llbracket p \bullet q \rrbracket x \equiv \text{Just } y$, we want to prove there exists z such that $\llbracket q \rrbracket x \equiv \text{Just } z$ and $\llbracket p \rrbracket z \equiv \text{Just } y$. Let σ be the result of *unify* $p_{\text{del}} q_{\text{ins}}$, witnessing $p \bullet q$; let γ be the result of *unify* $(\sigma q_{\text{del}}) x$, witnessing the application. Take $z = (\gamma \circ \sigma) q_{\text{ins}}$, and let us prove $\gamma \circ \sigma$ unifies p_{del} and z .

$$\begin{aligned}
 & (\gamma \circ \sigma) p_{\text{del}} \equiv (\gamma \circ \sigma) z && \{\text{z has no variables}\} \\
 \iff & (\gamma \circ \sigma) p_{\text{del}} \equiv z && \{\text{definition of } z\} \\
 \iff & (\gamma (\sigma p_{\text{del}})) \equiv \gamma (\sigma q_{\text{ins}}) && \{\text{hypothesis}\} \\
 \iff & \sigma p_{\text{del}} \equiv \sigma q_{\text{ins}}
 \end{aligned}$$

Hence, p can be applied to z , resulting in $(\gamma \circ \sigma) p_{\text{ins}}$, which is equal to y (hyp).

only if. Assuming there exists z such that $\llbracket q \rrbracket x \equiv \text{Just } z$ and $\llbracket p \rrbracket z \equiv \text{Just } y$, we want to prove that $\llbracket p \bullet q \rrbracket x \equiv \text{Just } y$. Let α be such that $\alpha q_{\text{del}} \equiv x$, hence, $z \equiv \alpha q_{\text{ins}}$; Let β be such that $\beta p_{\text{del}} \equiv z$, hence $y \equiv \beta p_{\text{ins}}$.

- a) First we prove that $p \bullet q$ is defined, that is, there exists σ' that unifies q_{ins} and p_{del} . Recall α and β have disjoint variables because we assume p and q have a disjoint set of names. Let $\sigma' = \alpha \cup \beta$, which corresponds to $\alpha \circ \beta$ or $\beta \circ \alpha$ because they have disjoint sets of names.

$$\sigma' q_{\text{ins}} \equiv \alpha q_{\text{ins}} \equiv z \equiv \beta p_{\text{del}} \equiv \sigma' p_{\text{del}}$$

Since σ' unifies q_{ins} and p_{del} , let σ be their *most general unifier*. Then, $\sigma' \equiv \gamma \circ \sigma$ for some γ and $p \bullet q \equiv \text{Chg}(\sigma q_{\text{del}})(\sigma p_{\text{ins}})$.

- b) Next we prove $\llbracket p \bullet q \rrbracket x \equiv \text{Just } y$. First we prove σq_{del} unifies with x .

$$\begin{aligned}
 & x \equiv \beta q_{\text{del}} && \{\text{Disj. supports; Def. } \sigma'\} \\
 \iff & x \equiv \gamma (\sigma q_{\text{del}}) && \{x \text{ has no variables}\} \\
 \iff & \gamma x \equiv \gamma (\sigma q_{\text{del}})
 \end{aligned}$$

Hence, $\llbracket p \bullet q \rrbracket x$ evaluates to $\gamma (\sigma p_{\text{ins}})$. Proving it coincides with y is a straightforward calculation:

$$\begin{aligned}
 & \gamma (\sigma p_{\text{ins}}) \equiv y && \{\text{Def. } y\} \\
 \iff & \gamma (\sigma p_{\text{ins}}) \equiv \alpha p_{\text{ins}} && \{\text{Disj. supports; Def. } \sigma'\} \\
 \iff & \gamma (\sigma p_{\text{ins}}) \equiv \gamma (\sigma p_{\text{ins}})
 \end{aligned}$$

□

Once we have established that composition is correct with respect to application, we would like to ensure composition is associative. But first we need to specify what we mean by *equal* changes. We will consider an extensional equality over changes. Two changes are said to be equivalent if and only if they are indistinguishable through their application semantics.

Definition 5.1.1 (Change Equivalent). Two changes p and q are said to be equivalent, denoted $p \approx q$, if and only if $\forall x . \llbracket p \rrbracket x \equiv \llbracket q \rrbracket x$

Lemma 5.1.2 (Definability of Composition). *Let p, q and r be aptly typed changes, then, $(p \bullet q) \bullet r$ is defined if and only if $p \bullet (q \bullet r)$ is defined.*

Proof. if. Assuming $(p \bullet q) \bullet r$ is defined, Let σ and θ be such that $\sigma p_{\text{del}} \equiv \sigma q_{\text{ins}}$ and $\theta (\sigma q_{\text{del}}) \equiv \theta r_{\text{ins}}$. We must prove that (a) r_{ins} unifies with q_{del} through some substitution θ' and (b) $\sigma' q_{\text{ins}}$ unifies with p_{del} . Take $\theta' = \theta \circ \sigma$, then:

$$\begin{aligned} & (\theta \circ \sigma) r_{\text{ins}} \equiv (\theta \circ \sigma) q_{\text{del}} & \{ \text{support } \sigma \cap \text{vars } r \equiv \emptyset \} \\ \iff & \theta r_{\text{ins}} \equiv (\theta \circ \sigma) q_{\text{del}} \end{aligned}$$

Let ζ be the idempotent *most general unifier* of r_{ins} and q_{del} , it follows that $\theta' = \gamma \circ \zeta$ for some γ . Consequently, $q \bullet r = \text{Chg } (\zeta r_{\text{del}}) (\zeta q_{\text{ins}})$.

Now, we must construct σ' to unify p_{del} and ζq_{ins} , which enables the construction of $p \bullet (q \bullet r)$. Let $\sigma' = \theta \circ \sigma$ and reduce it to one of our assumptions:

$$\begin{aligned} & \theta (\sigma p_{\text{del}}) \equiv \theta (\sigma (\zeta q_{\text{ins}})) & \{ \theta \circ \sigma \equiv \gamma \circ \zeta \} \\ \iff & \theta (\sigma p_{\text{del}}) \equiv \gamma (\zeta (\zeta q_{\text{ins}})) & \{ \zeta \text{ idempotent} \} \\ \iff & \theta (\sigma p_{\text{del}}) \equiv \gamma (\zeta q_{\text{ins}}) & \{ \theta \circ \sigma \equiv \gamma \circ \zeta \} \\ \iff & \theta (\sigma p_{\text{del}}) \equiv \theta (\sigma q_{\text{ins}}) \\ \iff & \sigma p_{\text{del}} \equiv \sigma q_{\text{ins}} \end{aligned}$$

only if. Analogous. □

Lemma 5.1.3 (Associativity of Composition). *Let p, q and r be aptly typed changes such that $(p \bullet q) \bullet r$ is defined, then $(p \bullet q) \bullet r \approx p \bullet (q \bullet r)$.*

Proof. Straightforward application of Lemma 5.1.2 and Lemma 5.1.1. □

Lemma 5.1.4 (Identity of Composition). *Let p be a change, then $\epsilon = \text{Chg } \#_x \#_x$ is the identity of composition. That is, $p \bullet \epsilon \approx p \approx \epsilon \bullet p$.*

Proof. Trivial; ϵ unifies with all possible terms. □

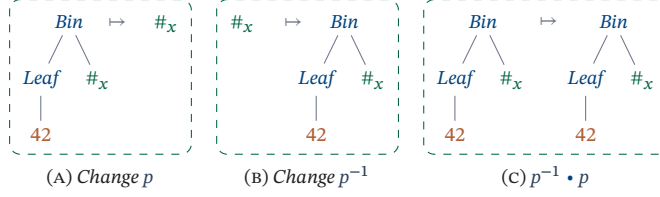


FIGURE 5.8: Example of a change, its inverse and their composition.

Lemmas 5.1.3 and 5.1.4 establish a partial monoid structure for Chg and \cdot under extensional change equality, \approx . It is not trivial to squeeze more structure out of this change representation, as we shall discuss next.

LOOSE ENDS. The first thing that comes to mind is the definition of the inverse of a change. Since changes are well-scoped, that is, $\text{vars } p_{\text{del}} \equiv \text{vars } p_{\text{ins}}$ for any change p , defining the inverse of a change p , denoted p^{-1} , is trivial:

$$\begin{aligned} \cdot^{-1} &:: \text{Chg } \kappa \text{ fam at} \rightarrow \text{Chg } \kappa \text{ fam at} \\ p^{-1} &= \text{Chg } p_{\text{ins}} p_{\text{del}} \end{aligned}$$

Naturally, then, we would expect that $p \cdot p^{-1} \approx \epsilon$, but that is not the case. The domain of ϵ is the entire set of trees, but the domain of $p \cdot p^{-1}$ is generally strictly smaller. Consequently, we can easily find a tree t such that $\llbracket \epsilon \rrbracket t \equiv \text{Just } t$ but $\llbracket p \cdot p^{-1} \rrbracket t \equiv \text{Nothing}$. Take, for example, the change shown in Figure 5.8.

The problem with inverses above stems from $p \cdot p^{-1}$ being *less general* than the identity, since it has a smaller domain. In other words, $p \cdot p^{-1}$ works on a subset of the domain of ϵ , but it performs the same action as ϵ for the elements it is defined. It is natural then to attempt to talk about changes modulo their domain. We could think of stating $p \leq q$ whenever $\llbracket p \rrbracket \subseteq \llbracket q \rrbracket$. That is, when p and q are the same except that the domain of q is larger. This \leq is known as the usual *extension order* [93], and when instantiated for our particular case, yields the definition below.

Definition 5.1.2 (Extension Order). Let p and q be two aptly typed changes; we say that q is an extension of p , denoted $p \leq q$, if and only if $\forall x \in \text{dom } p . \llbracket p \rrbracket x \equiv \llbracket q \rrbracket x$. In other words, $p \leq q$ when q coincides with p when restricted to p 's domain.

This gives us a preorder on changes and it is the case that $p \cdot p^{-1} \leq \epsilon$ and $p^{-1} \cdot p \leq \epsilon$. Attempting to identify $p \cdot p^{-1}$ as somehow equivalent to ϵ using \leq will not work, however. We could think of defining a notion of *approximate changes*, denoted $p \sim q$, by whether p and q are comparable under \leq . This would not yield an equivalence relation since \sim is not transitive, as illustrated in Figure 5.9). Moreover, the extension order cannot be used

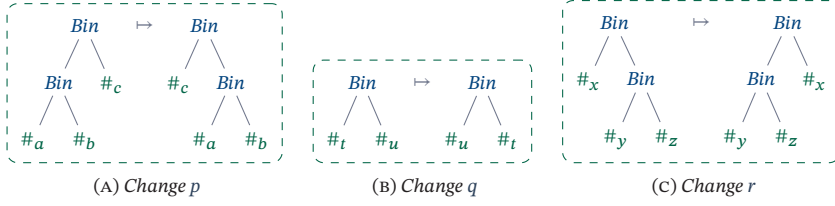


FIGURE 5.9: Three changes such that $p \sim q$ (because $p \leq q$) and $q \sim r$ (because $r \leq q$). Yet, $p \not\sim r$ since its not the case that $r \leq p$ or $p \leq r$ holds.

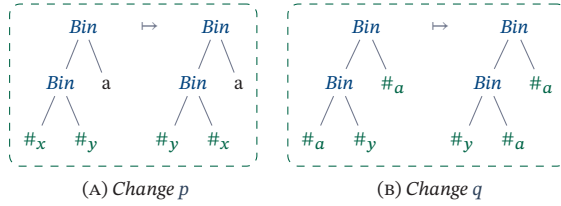


FIGURE 5.10: Two changes that could be used to transform the same element x but are not comparable under the extension order (\leq).

to define the *best* change between two elements x and y . Take x to be `Bin (Bin a b) a` and y to be `Bin (Bin b a) a`, two uncomparable changes are shown in Figure 5.10.

This short discussion does not mean that there is *no* suitable way to compare the changes in Figure 5.10 or to define \sim in such a way that the changes in Figure 5.9 be considered equivalent. It does mean, however, that comparing the domain of changes is a weak definition and a robust definition will probably be significantly more involved.

5.1.4 COMPUTING CHANGES

Having seen how `Chg` has the basic properties we would expect, we move on to computing them. In this section we define the generic counterpart to the `chgTree` function (Section 5.1.1). Recall that the differencing algorithm starts by computing the *sharing map* of its source s and destination d , which enable us to efficiently decide if a given tree x is a subtree of s and d . Later, we use this sharing map and *extract* the deletion and insertion contexts, according to some extraction mode, which ensure we will produce well-scoped changes (Figure 5.4).

data `ExtractionMode = NoNested | ProperShare | Patience`

The *sharing map* is central to the efficiency of the differencing algorithm, but it marks subtrees for sharing regardless of underlying semantics, which can be a problem when the trees in question represent complex structures such as abstract syntax trees. We must be careful not to *overshare* trees. Imagine a local variable declaration `int x = 0;` inside an arbitrary function. This declaration should *not* be shared with another syntactically equal declaration in another function. A careful analysis of what can and cannot be shared would require domain-specific knowledge of the programming language in question. Nevertheless, we can impose different restrictions that make it *unlikely* that values will be shared across scope boundaries. A simple and effective such measure is not sharing subtrees with height strictly less than one (or a configurable parameter). This keeps constants and most variable declarations from being shared, effectively avoiding the issue.

5.1.4.1 WHICH COMMON SUBTREE, GENERICALLY

Similarly to the example from Section 5.1.1, the first thing we must do is to annotate our trees with hashes at every point. The *Holes* datatype from `generics-simplistic` also supports annotations. Unlike the concrete example, however, we will also keep the height of each tree to enable us to easily forbid sharing trees smaller than a certain height. The *PrepFix* datatype, defined below, serves the same purpose as the simpler *TreeH*, from our concrete example.

```
data PrepData a = PrepData {getDigest :: Digest, getHeight :: Int}
type PrepFix κ fam = SFixAnn κ fam PrepData
```

The *decorate* function can be written with the help of synthesized attributes (Section 3.2.3.1). The homonym *synthesize* function from `generics-simplistic` serves this very purpose. We omit the algebra passed to *synthesize* but invite the interested reader to check *Data.HDiff.Diff.Preprocess* in the source (Appendix A).

```
decorate :: (All Digestible κ) ⇒ SFix κ fam at → PrepFix κ fam at
decorate = synthesize ...
```

The algebra used by *decorate*, above, computes a hash at each constructor of the tree. The hashes are computed from a unique identifier per constructor and a concatenation of the hashes of the subtrees. The hash of the root in Figure 5.6, for example, is computed with a call to *hash* (*concat* [“*Main.Tree.Bin*”, “*310dac*”, “*4a32bd*”]). This ensures that hashes uniquely identify a subtree modulo hash collisions.

After preprocessing the input trees we traverse them and insert every hash we see in a hash map from hashes to integers. These integers count how many times we have seen a tree, indicating the arity of a subtree. Shared subtrees occur with arity of at least

two: once in the deletion context and once in the insertion context. The underlying datastructure is a *Int64*-indexed trie [19].

```
type Arity = Int
buildArityMap :: PrepFix a  $\kappa$  fam ix  $\rightarrow$  Trie Arity
```

A call to *buildArityMap* with the annotated tree shown in Figure 5.6, for example, would yield the map *fromList* [(*“0f42ab”*, 1), (*“310dac”*, 1), (*“0021ab”*, 2), ...].

After processing the *arity* maps for both the source tree and destination tree, we construct the *sharing* map, which consists in the intersection of the arity maps and a final pass adding a unique identifier to every key. We also keep track of how many metavariables were assigned, so we can always allocate fresh names without having to inspect the whole map again. This is just a technical consequence of working with binders explicitly.

```
type MetavarAndArity = MAA {getMetavar :: Int, getArity :: Arity}
buildSharingMap :: PrepFix a  $\kappa$  fam ix  $\rightarrow$  PrepFix a  $\kappa$  fam ix
                   $\rightarrow$  (Int, Trie MetavarAndArity)
buildSharingMap x y = T.mapAccum ( $\lambda$  i ar  $\rightarrow$  (i + 1, MAA i ar)) 0
                      $ T.zipWith (+) (buildArityMap x) (buildArityMap y)
```

The final *wcs* *s* *d* is straightforward: we preprocess the trees with their hash and height; then compute their sharing map, which is used to lookup the common subtrees. Yet, the whole point of preprocessing the trees was to avoid the unnecessary recomputation of their hashes. Consequently, we are better off carrying these preprocessed trees everywhere through the computation of changes. The final *wcs* function will have its type slightly adjusted and is defined below.

```
wcs :: (All Digestible  $\kappa$ )  $\Rightarrow$  PrepFix  $\kappa$  fam at  $\rightarrow$  PrepFix  $\kappa$  fam at
       $\rightarrow$  PrepFix  $\kappa$  fam at  $\rightarrow$  Maybe Int
wcs s d = let m = buildSharingMap s d
          in famp getMetavar  $\circ$  flip T.lookup m  $\circ$  getDigest  $\circ$  getAnnot
```

Let $f = \text{wcs } s \ d$ for some *s* and *d*. Computing *f* itself is linear and takes $\mathcal{O}(n + m)$ time, where *n* and *m* are the number of constructors in *s* and *d*. A call to *f* *x* for some *x*, however, is answered in $\mathcal{O}(1)$ due to the bounded depth of the trie.

We chose to use a cryptographic hash function [69] and ignore the remote possibility of hash collisions. Although it would not be hard to detect these collisions whilst computing the arity map, doing so would incur a performance penalty. Checking for collisions would require us to store the path to the tree together with its associated hash instead of only storing the hash. Then, on every insertion we could check that the inserted tree

matches with the tree already in the map. Had we opted for a non-cryptographic hash, which are much faster to compute than cryptographic hash functions, we would have had to employ the collision detection mechanism above. It is plausible that this would cost more time than computing the cryptographic hash at once. We did not test this, however.

5.1.4.2 CONTEXT EXTRACTION

After computing the set of common subtrees, we must decide which of those subtrees should be shared. Shared subtrees are abstracted by a metavariable in every location they would occur at in the deletion and insertion contexts.

Recall that we chose to never share subtrees with height smaller than a given parameter. Our choice is very pragmatic in the sense that we can preprocess the necessary information and it effectively avoids most of the oversharing without involving domain specific knowledge. The *CanShare* below is a synonym for a predicate over trees used to decide whether we can share a given tree or not.

type *CanShare* κ *fam* = $\forall ix \cdot \text{PrepFix } \kappa \text{ fam } ix \rightarrow \text{Bool}$

The *extract* function takes an *ExtractionMode*, a sharing map and a *CanShare* predicate and two preprocessed fixpoints to extract contexts from. The reason we receive two trees at the same time and produce two contexts is because modes like *NoNested* perform some cleanup that depends on global information.

extract :: *ExtractionMode* \rightarrow *CanShare* κ *fam* \rightarrow *IsSharedMap*
 $\rightarrow (\text{PrepFix } \kappa \text{ fam} :: \text{PrepFix } \kappa \text{ fam}) \text{ at} \rightarrow \text{Chg } \kappa \text{ fam at}$

To some extent, we could compare context extraction to the translation of tree mappings into edit-scripts: our tree matching is encoded in *wcs* and instead of computing an edit-script, we compute terms with metavariables. Classical algorithms are focused in computing the *least cost* edit-script from a given tree mapping. In our case, the notion of *least cost* hardly makes sense – besides not having defined a cost semantics to our changes, we are interested in those that merge better which might not necessarily be those that insert and delete the least amount of constructors. Consequently, there is a lot of freedom in defining our context extraction techniques. We will look at three particular examples next, but we discuss other possibilities later (Section 5.4).

EXTRACTING WITH *NoNested*. Extracting contexts with the *NoNested* mode happens in two passes. We first extract the contexts naively, then make a second pass removing the variables that appear exclusively in the insertion. To keep the extraction algorithm

linear it is important to *not* forget which common subtrees have been substituted on the first pass. Hence, we create a context that contains metavariables and their associated tree.

```

noNested1 :: CanShare  $\kappa$  fam  $\rightarrow$  Trie MetavarAndArity  $\rightarrow$  PrepFix  $\kappa$  fam at
            $\rightarrow$  Holes  $\kappa$  fam (Const Int  $:\ast$ : PrepFix a  $\kappa$  fam) at
noNested1 h sm x@(PrimAnn _ xi) = Prim xi
noNested1 h sm x@(SFixAnn ann xi)
    = if h x then maybe recurse (mkHole x) $ lookup (getDigest ann) sm
      else recurse
where recurse    = Roll (repMap (noNested1 h sm) xi)
      mkHole x v = Hole (Const (getMetavar v)  $:\ast$ : x)

```

The second pass maps over the holes in the output from the first pass and decides whether to transform the *Const Int* into a *Metavar κ fam* or whether to forget this was a potential shared tree and keep the tree instead. We will omit the implementation of the second pass. It consists in a straightforward traversal of the output of *noNested1*. We direct the interested reader to check *Data.HDiff.Diff.Modes* in the source code for more details (Appendix A).

EXTRACTING WITH *Patience*. The *Patience* extraction can be done in a single pass. Unlike *noNested1* above, instead of simply looking a hash up in the sharing map, it will further check that the given hash occurs with arity two – indicating the tree in question occurs once in the source tree and once in the destination. This completely bypasses the issue with *NoNested* producing insertion contexts with undefined variables and requires no further processing. The reason for it is that the variables produced will appear with the same arity as the trees they abstract, twice in this case: once in the deletion context and once in the insertion context.

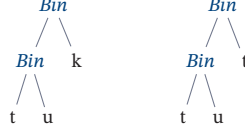
```

patience :: CanShare  $\kappa$  fam  $\rightarrow$  Trie MetavarAndArity  $\rightarrow$  PrepFix a  $\kappa$  fam at
            $\rightarrow$  Holes  $\kappa$  fam (Metavar  $\kappa$  fam) at
patience h sm x@(PrimAnn _ xi) = Prim xi
patience h sm x@(SFixAnn ann xi)
    = if h x then maybe recurse (mkHole x) $ lookup (getDigest ann) sm
      else recurse
where recurse    = Roll (repMap (patience h sm) xi)
      mkHole x v | getArity v  $\equiv$  2 = Hole ( $\#_{\text{(getMetavar v)}}^{\text{fam}}$ )
      | otherwise    = $fixToHoles x

```

EXTRACTING WITH *ProperShares*. The *ProperShares* method prefers sharing smaller subtrees more times instead of bigger subtrees, which might shadow nested commonly occurring subtrees (Figure 5.3).

Given a source s and a destination d , we say that a tree x is a *proper-share* between s and d whenever no subtree of x occurs in s and d with arity greater than that of x . In other words, x is a proper-share if and only if all of its subtrees occur only as subtrees of other occurrences of x . For the two trees below, u is a proper-share but $\text{Bin } t \ u$ is not: t occurs once *outside* $\text{Bin } t \ u$.



Extracting contexts with under the *ProperShare* mode consists in annotating the source and destination trees with a boolean indicating whether or not they are a proper share, then proceeding just like *Patience*, but instead of checking that the arity must be two, we check that the tree is classified as a *proper-share*. It is important to use annotated fixpoints to maintain performance, but the code is very similar to the previous two methods and, hence, omitted.

THE *chg* FUNCTION. Finally, the generic *chg* function receives a source and destination trees, s and d , and computes a change that encodes the information necessary to transform the source into the destination according to some extraction mode *extMode*. In our prototype, the extraction mode comes from a command line option.

```

chg :: (All Digestible κ) ⇒ SFix κ fam at → SFix κ fam at → Patch κ fam at
chg x y = let dx      = decorate x
          dy      = decorate y
          (_, sh) = buildSharingMap opts dx dy
          in extract extMode canShare (dx :* dy)
where
  canShare t = 1 < treeHeight (getConst (getAnn t))
  
```

5.2 THE TYPE OF PATCHES

Up until now we have seen how *changes* consisting of a deletion and an insertion context are a suitable representation for encoding transformations between trees. In fact, changes are very similar to *tree matchings* (Section 2.1.2) but with fewer restrictions. From a synchronization point of view, however, these *changes* are very difficult to merge. They do not explicitly encode enough information for that.

Synchronizing changes requires us to understand which constructors in the deletion context are, in fact, just being copied over in the insertion context. Taking Figure 5.11,

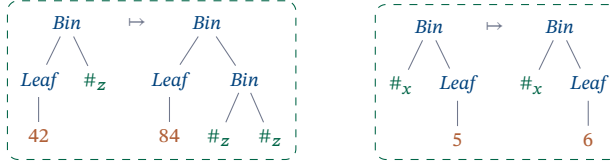


FIGURE 5.11: Example of disjoint changes. Each change is delimited by a dashed box. The leftmost change modifies the left child and duplicates the right child without changing its content. The rightmost change operates solely on the right child.

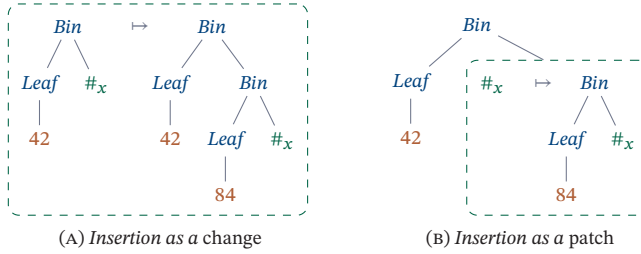


FIGURE 5.12: A change with redundant information on the left and its minimal representation on the right, with an evident spine.

where one change operates exclusively on the right child of a binary tree whereas the other alters the left child and duplicates the right child in-place. These changes are clearly *disjoint*, since they modify the content of different subtrees of the source. Consequently it should be possible to be automatically synchronize them. To recognize them as *disjoint* changes, though, will require more information than what is provided by *Chg*.

Observing the definition of *Chg* reveals that the deletion context might *delete* many constructors that are being inserted, in the same place, by the insertion context. The changes from Figure 5.11, for example, conceal the fact that the *Bin* at the root of the source tree is, in fact, being copied in both changes. Following the *stdiff* nomenclature, the *Bin* at the root of both changes in Figure 5.11 should be placed in the *spine* of the patch. That is, it is copied over from source to destination but it leads to changes further down the tree.

A *patch*, then, captures the idea of many individual changes operating over separate parts of the source tree. It consists in a spine that leads to changes in its leaves, and is defined by the type *Patch* below.

type *Patch* κ *fam* = *Holes* κ *fam* (*Chg* κ *fam*)

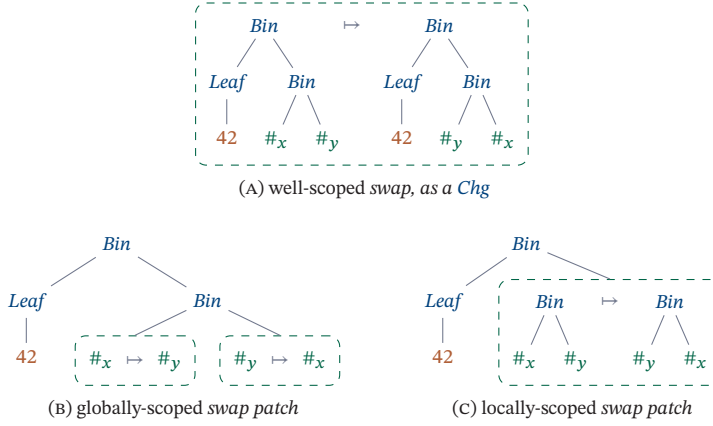


FIGURE 5.13: A change that swaps some elements; naive anti-unification of the deletion and insertion context breaking scoping; and finally the patch with minimal changes.

Figure 5.12 illustrates the difference between patches and changes. In Figure 5.12(A) we see *Bin* (*Leaf* 42) being repeated in both contexts – whereas in Figure 5.12(B) it has been placed in the spine and consequently, is clearly identified as a copy.

Patches are computed from changes by extracting common constructors from the deletion and insertion contexts into the spine. In other words, we would like to push the changes down towards the leaves of the tree. There are two different ways for doing so, illustrated by Figure 5.13. On one hand we can consider the patch metavariables to be *globally-scoped*, yielding structurally minimal changes, Figure 5.13(B). On the other hand, we could strive for *locally-scoped*, where each change might still contain repeated constructors as long as they are necessary to ensure the change is *closed*, as in Figure 5.13(C). The first option, *globally-scoped* patches, is very easy to compute. All we have to do is to compute the anti-unification of the insertion and deletion context.

$globallyScopedPatch :: Chg \text{ ki codes at} \rightarrow Patch_{PE} \text{ ki codes at}$
 $globallyScopedPatch (Chg \ d \ i) = holesMap (uncurry' \ Chg) (lgg \ d \ i)$

Globally-scoped patches are easy to compute but contribute little information from a synchronization point of view. To an extent, it makes merging even harder. Take Figure 5.14, where a globally scoped patch is produced from a change. It is harder to understand that the (42:) is being deleted by looking at the globally-scoped patch than by looking at the change. This is because the first (:) constructor is considered to be in the spine by the naive anti-unification algorithm, which proceeds top-down. A bottom-up approach is also unpractical, as we would have to decide which leaves to pair together and it would suffer similar issues for data inserted on the tail of linearly structured data.

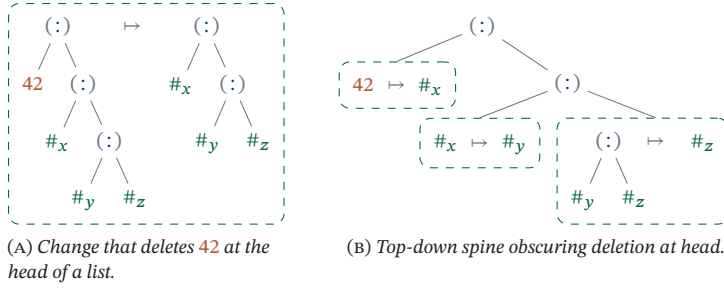


FIGURE 5.14: Globally-scoped patches resulting in misalignment of metavariables due to deletions in the head of linearly-structured data.

Locally-scoped patches imply that changes might still contain repeated constructors in the root of their deletion and insertion contexts – hence they will not be structurally minimal. Although more involved to compute, they give us a chance to address insertions and deletions of constructors before we end up misaligning copies.

Independent of global or local scoping, ignoring the information about the spine yields a forgetful functor from patches back into changes, named *chgDistr*. Its definition is straightforward thanks to the free monad structure of *Holes*, which gives us the necessary monadic multiplication. We must take care that *chgDistr* will not capture variables, that is, all metavariables must have already been properly α -converted. We cannot enforce this invariant directly in the *chgDistr* function for performance reasons, consequently, we must manually ensure that all scopes contain disjoint sets of names and therefore can be safely distributed whenever applicable. This is a usual difficulty when handling objects with binders, in general.

$$\begin{aligned}
 \text{holesMap} &:: (\forall x . \varphi x \rightarrow \psi x) \Rightarrow \text{Holes } \kappa \text{ fam } \varphi \text{ at} \rightarrow \text{Holes } \kappa \text{ fam } \psi \text{ at} \\
 \text{holesJoin} &:: \text{Holes } \kappa \text{ fam } (\text{Holes } \kappa \text{ fam}) \text{ at} \rightarrow \text{Holes } \kappa \text{ fam } \text{ at} \\
 \text{chgDistr} &:: \text{Patch } \text{ki codes at} \rightarrow \text{Chg } \text{ki codes at} \\
 \text{chgDistr } p &= \text{Chg } (\text{holesJoin } (\text{holesMap } \cdot_{\text{del}} p)) (\text{holesJoin } (\text{holesMap } \cdot_{\text{ins}} p))
 \end{aligned}$$

The application semantics of *Patch* is independent of the scope choices, and is easily defined in terms of *chgApply*. First we compute a global change that corresponds to the patch in question, then use *chgApply*. The *apply* function below works for locally and globally scoped patches, as long as we care that the precondition for *chgDistr* is maintained.

$$\begin{aligned}
 \text{apply} &:: (\text{All Eq } \kappa) \Rightarrow \text{Patch } \kappa \text{ fam at} \rightarrow \text{SFix } \kappa \text{ fam at} \rightarrow \text{Maybe } (\text{SFix } \kappa \text{ fam at}) \\
 \text{apply } p &= \text{chgApply } (\text{chgDistr } p)
 \end{aligned}$$

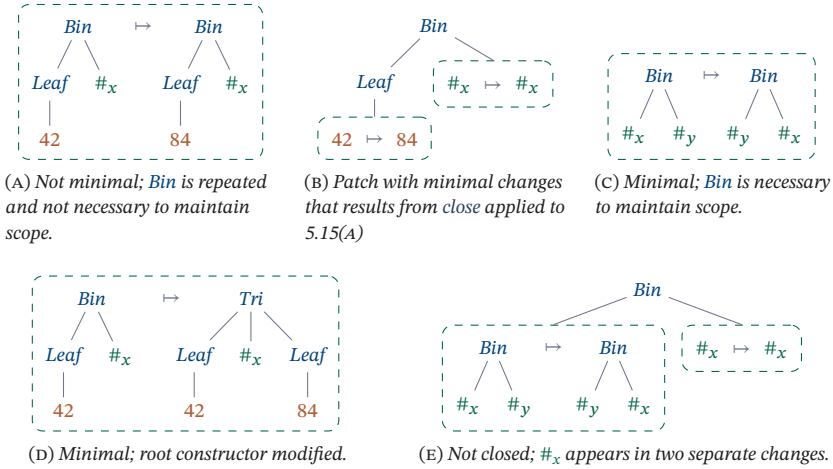


FIGURE 5.15: Some non-minimal-closed and minimal-closed changes examples.

Overall, we find ourselves in a dilemma. On the one hand we have *globally-scoped* patches, which have larger spines but can produce results that are difficult to understand and synchronize, as in Figure 5.14. On the other hand, *locally-scoped* patches are more involved to compute, as we will study next, in Section 5.2.1, but they forbid misalignments and also enable us to process small changes independently of one another in the tree. This is particularly important for being able to develop an industrial synchronizer at some point, as it keeps *conflicts* small and isolated.

We propose that the actual solution will consist in using a combination of both local and global scoping. First we will produce a locally-scoped patch, which forbids situations as in Figure 5.14. This patch will consist in an *outer* spine leading to closed locally-scoped changes. This gives us an opportunity to identifying deletions and insertions that could cause copies to be misaligned, essentially producing a globally-scoped *alignment* inside each of those changes. Alignments will be discussed in more detail shortly (Section 5.2.3).

5.2.1 COMPUTING CLOSURES

Computing locally-scoped patches consists of first computing the largest possible spine, like we did with globally-scoped patches, then enlarging the resulting changes until they are well-scoped and closed. Figure 5.13 illustrates this process. Computing the closure of Figure 5.13(A) starts with Figure 5.13(B), then *enlarging* the changes so that they contain the *Bin* constructor, which fixes their scope (resulting in Figure 5.13(C)).

We say a change is closed when it has no free metavariables and, additionally, its metavariables occur nowhere else. The changes produced by the *chg* function are closed, for example, but they might not be as small as they could be. We say a change is *minimal* when the root constructors in its deletion and insertion context are either different or necessary to maintain scope. Figure 5.15 illustrates different combinations of *closed* and *minimal* changes. The intuition behind *minimal-closed* changes is that two such changes should not interfere with one another.

Producing locally-scoped minimal-closed changes can be difficult under arbitrary renamings. Take Figure 5.15(E), one could argue that if the occurrences of $\#_x$ in each individual change are, in fact, different, then the changes are minimal-closed. To avoid this, we always start from a large well-scoped change produced with *chg*. Consequently, we know that every occurrence of $\#_x$ refers to *the same* tree in the source of the patch. This is another technicality of dealing with names explicitly and provides good reason to enforce that names are always different, even when occurring in separate scopes.

In general, then, we can only know that a change is in fact closed if we know how many times each variable is used globally. Say a variable $\#_z$ is used $n + m$ times in total within a change c , and it has n and m occurrences in the deletion and insertion contexts of c , respectively. Then $\#_z$ does not occur anywhere else but within c , in other words, $\#_z$ is *local* to c . If all variables of c are *local* to c with respect to some global scope, we say c is closed. Given a multiset of variables for the global scope, we can define *isClosedChg* in Haskell as:

```
isClosedChg :: Map Int Arity → Chg κ fam at → Bool
isClosedChg global (Chg d i) = isClosed global (vars d) (vars i)
  where isClosed global ds us = unionWith (+) ds us `isSubmapOf` global
```

The *close* function, shown in Figure 5.16, is responsible for pushing constructors through the least general generalization until they represent minimal-closed changes. It calls an auxiliary version that receives the global scope and keeps track of the variables it has seen so far. The worst case scenario happens when we need *all* constructors of the spine to close the change, in which case, *close c = Hole c*; yet, if we pass a non-well-scoped change to *close*, it cannot produce a result and throws an error instead.

Efficiently computing closures requires us to keep track of the variables that have been declared and used in a change – that is, we have seen occurrences in the deletion and insertion context respectively. Recomputing these multisets would result in a slower algorithm. The *annWithVars* function below computes the variables that occur in two contexts and annotates a change with them:

```
data WithVars x at = WithVars { decls , uses :: Map Int Arity , body :: x at }
withVars :: (HolesMV κ fam ::*: HolesMV κ fam) at → WithVars (Chg κ fam) at
withVars (d ::*: i) = WithVars (vars d) (vars i) (Chg d i)
```

```

close :: Chg  $\kappa$  fam at  $\rightarrow$  Holes  $\kappa$  fam (Chg  $\kappa$  fam) at
close c@(Chg d i) = case closeAux (chgVars c) (holesMap withVars (lgg d i)) of
  InL _  $\rightarrow$  error "invariant failure: c was not well-scoped."
  InR b  $\rightarrow$  holesMap body b

closeAux :: M.Map Int Arity  $\rightarrow$  Holes  $\kappa$  fam (WithVars (Chg  $\kappa$  fam)) at
 $\rightarrow$  Sum (WithVars (Chg  $\kappa$  fam)) (Holes  $\kappa$  fam (WithVars (Chg  $\kappa$  fam))) at
closeAux _ (Prim x) = InR (Prim x)
closeAux gl (Hole cv) = if isClosed gl cv then InR (Hole cv) else InL cv
closeAux gl (Roll x) =
  let aux = repMap (closeAux gl) x
  in case repMapM fromInR aux of
    Just res  $\rightarrow$  InR (Roll res)
    Nothing  $\rightarrow$  let res = chgVarsDistr (Roll (repMap (either' Hole id) aux))
      in if isClosed gl res then InR (Hole res) else InL res
where
  fromInR :: Sum f g x  $\rightarrow$  Maybe (g x)

```

FIGURE 5.16: Complete generic definition of *close* and *closeAux*.

The *chgVarsDistr* is the engine of the *close* function. It distributes a spine over a change, similar to *chgDistr*, but here we care to maintain the explicit variable annotations correctly.

```

chgVarsDistr :: Holes  $\kappa$  fam (WithVars (Chg  $\kappa$  fam)) at  $\rightarrow$  WithVars (Chg  $\kappa$  fam) at
chgVarsDistr rs = let us = map (exElim uses) (getHoles rs)
  ds = map (exElim decls) (getHoles rs)
  in WithVars (unionsWith (+) ds) (unionsWith (+) us)
  (chgDistr (repMap body rs))

```

The *closeAux* function, Figure 5.16, receives a spine with leaves of type *WithVars ...* and attempts to *enlarge* them as necessary. If it is not possible to close the current spine, we return a *InL ...* equivalent to pushing all the constructors of the spine down the deletion and insertion contexts.

5.2.2 THE *diff* FUNCTION

Equipped with the ability to produce changes and minimize them, we move on to defining the *diff* function. As usual, it receives source and destination trees, *s* and *d*, and computes a patch that encodes the information necessary to transform the source into

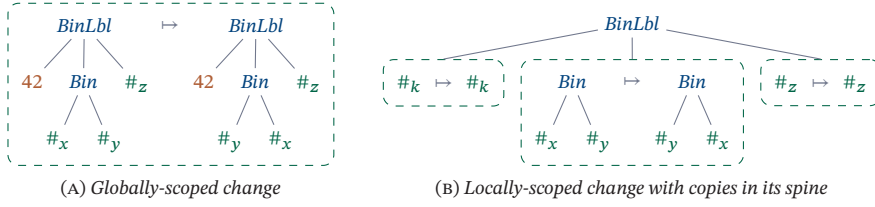


FIGURE 5.17: A Globally-scoped change and the result of applying it to *cpyPrimsOnSpine* \circ *close*, producing a patch with locally scoped changes and copies in its spine.

the destination. The extraction of the contexts yields a *Chg*, which is finally translated to a locally-scoped *Patch* by identifying the largest possible spine, with *close*.

```

diff :: (All Digestible  $\kappa$ )  $\Rightarrow$  SFix  $\kappa$  fam at  $\rightarrow$  SFix  $\kappa$  fam at  $\rightarrow$  Patch  $\kappa$  fam at
diff x y = let dx      = preprocess x
            dy      = preprocess y
            (i, sh)   = buildSharingMap opts dx dy
            (del :*: ins) = extract extMode canShare (dx :*: dy)
            in cpyPrimsOnSpine i (close (Chg del ins))
where canShare t = 1 < treeHeight (getConst (getAnn t))

```

The *cpyPrimsOnSpine* function will issue copies for the opaque values that appear on the spine, as illustrated in Figure 5.17. There, the *42* does not get shared for its height is smaller than 1 but since it occurs in the same location in the deletion and insertion context it can be identified as a copy – which involves issuing a fresh metavariable, hence the parameter *i* in the code above.

5.2.3 ALIGNING CHANGES

As we have seen in the previous sections, locally-scoped changes can avoid misaligning changes (Figure 5.14), but they still do not help us in identifying the insertions and deletions. As it will turn out, identifying these insertions and deletions is crucial for synchronization. In this section we will define a datatype and an algorithm for representing and computing alignments, which make the backbone of synchronization. Untyped synchronizers, such as *harmony* [35], must employ schemas to identify insertions and deletions avoiding misalignments (Figure 5.14). In our case, the type information enables us to identify insertions and deletions naturally by ensuring that they delete one layer of a *recursive type* at a time, never altering the type of the value under scrutiny.

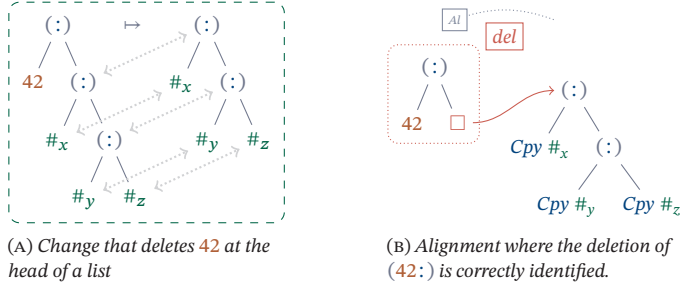


FIGURE 5.18: The change from Figure 5.14, with with an association of which nodes of the deletion and insertion contexts represent the same information, and an explicit representation of that information.

Take Figure 5.18(A), illustrating the change that motivated locally-scoped patches (Figure 5.14) in the first place. This time, however, arrows connect constructors that represent *the same information* in each respective context. This makes it clear that $(42 :)$ has no counterpart in the insertion context and, consequently, corresponds to a deletion. The *Chg* datatype by itself is insufficient to represent all this information. Therefore we need a new datatype for *alignments*, *Al*, and a function that translates a *Chg* into an *Al*. Computing and representing an alignment is, intuitively, the process of computing and representing this association between subtrees of the deletion and insertion contexts. The aligned version of Figure 5.18(A) is shown in Figure 5.18(B), where the *Al* border marks scoping for metavariables. The constructors that are paired up in the deletion and insertion are placed in a spine; those without a correspondent are flagged as deletions or insertions depending on which context they belong. Finally, *Cpy* $\#_{\square}$ is an abbreviation for *Chg* $\#_{\square} \#_{\square}$.

An aligned patch consists of a spine of copied constructors leading to a *well-scoped alignment*. This alignment, in turn, consists of a sequence of insertions, deletions or spines, which finally lead to a *Chg*. These *Chg* in the leaves of the alignment are globally-scoped with respect to the alignment they belong. We also add explicit information about copies and permutations to aid the synchronization engine later. Figure 5.19 illustrates a value of type *Patch* and its corresponding alignment, of type *PatchAl* defined below. Note how the scope from each change in Figure 5.19(A) is preserved in Figure 5.19(B), but the *Bin* on the left of the root can now be safely identified as a copy without losing information about the scope of $\#_x$.

type *PatchAl* κ fam = *Holes* κ fam (*Al* κ fam (*Chg* κ fam))

Computing the *alignment* for a change c consists in identifying what information in the deletion context correspond to *the same information* in the insertion context. The

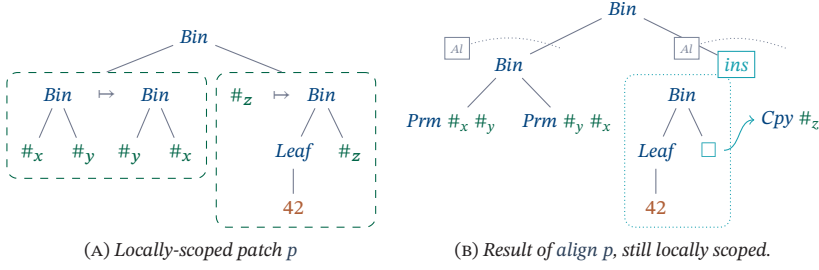


FIGURE 5.19: A patch p and its corresponding aligned version. The AI barrier marks the beginning of an alignment and delimits scopes; copies and permutations are marked explicitly and insertions and deletions indicate their continuation with \square .

bits and pieces in the deletion context that have no correspondent in the insertion context should be identified as deletions and vice-versa for the insertion context. In Figure 5.18(A), for example, the second $(:)$ in the deletion context represents the same information as the root $(:)$ in the insertion context.

We can recognize the deletion of $(42:)$ in Figure 5.18(B) structurally. All of its fields, except one recursive field, contains no metavariables. The one subtree which *does contain* metavariables is denoted the *focus* of the deletion (resp. insertion). We denote trees with no metavariables as *rigid* trees. A *rigid* tree has the guarantee that none of its subtrees is being copied, moved or modified. Consequently, *rigid* trees are being entirely deleted from the source or inserted at the destination of the change. If a constructor in the deletion (resp. insertion) context has all but one of its subtrees being *rigid*, it is only natural to consider this constructor to be part of the *deletion* (resp. *insertion*).

Since our patches are locally scoped, computing an aligned patch is done by mapping over the spine and aligning the individual changes. Aligning changes, in turn, consists in identifying whether the constructor at the head of the deletion (resp. insertion) context can be deleted (resp. inserted) then recursing on the focus of the deletion (resp. insertion). When the root of the deletion context and the root of the insertion context qualify for deletion and insertion, we check whether we can add them to a spine instead.

5.2.3.1 GENERIC ALIGNMENTS

We will be representing a deletion or insertion of a recursive *layer* by identifying the *position* where this modification must take place. Moreover, said position must be a recursive field of the constructor – that is, the deletion or insertion must not alter the type that our patch operates over. This is easy to identify since we follow a typed approach, where we always have access to type-level information.

In the remainder of this section we discuss the datatypes necessary to represent an aligned change, as illustrated in Figure 5.18(B), and how to compute said alignments from a $\text{Chg } \kappa \text{ fam } at$. The alignChg function, declared below, will receive a well-scoped change and compute an alignment.

$$\text{alignChg} :: \text{Chg } \kappa \text{ fam } at \rightarrow \text{Al } \kappa \text{ fam } (\text{Chg } \kappa \text{ fam}) at$$

Alignments encoded in the Al datatype are similar to its predecessor $\text{Al}\mu$ from `stdiff` (Section 4.1.2). They record insertions, deletions and spines over a fixpoint. Insertions and deletions are represented with *Zipper*s [45]. A zipper over a datatype t is the type of *one-hole-contexts* over t , where the hole indicates a focused position. We will use the zippers provided directly by the `generics-simplistic` library (Section 3.2.4.1). These zippers encode a *single* layer of a fixpoint at a time, for example, a zipper over the Bin constructor is either $\text{Bin } \square u$ or $\text{Bin } u \square$, indicating the focus is in either the left or the right subtree. It *does not* enable us specify a nested focus point, like in $\text{Bin } (\text{Bin } \square t) u$.

A value of type $\text{Zipper } c g h at$ is then equivalent to a constructor of type at with one of its recursive positions replaced by a value of type $h at$ and the other positions at' (recursive or not) carrying values of type $g at'$. The c above is a constraint that enables us to inform GHC about some properties of type at and is mostly a technicality.

An alignment $\text{Al } \kappa \text{ fam } f at$ represents a sequence of insertions and deletions interleaved with spines, copies and permutations which ultimately lead to *unclassified modifications*, which are typed according to the f parameter. Next, we will go through the six constructors of Al one by one. First we have deletions and insertions, which explicitly mention a zipper and one recursive field to continue the alignment.

data $\text{Al } \kappa \text{ fam } f at$ **where**

$$\text{Del} :: \text{Zipper } (\text{CompoundCnstr } \kappa \text{ fam } at) (\text{SFix } \kappa \text{ fam}) (\text{Al } \kappa \text{ fam } f) at \rightarrow \text{Al } \kappa \text{ fam } f at$$

$$\text{Ins} :: \text{Zipper } (\text{CompoundCnstr } \kappa \text{ fam } at) (\text{SFix } \kappa \text{ fam}) (\text{Al } \kappa \text{ fam } f) at \rightarrow \text{Al } \kappa \text{ fam } f at$$

The CompoundCnstr constraint above must be carried around to indicate we are aligning a type that belongs to the mutually recursive family and therefore has a generic representation – again, just a Haskell technicality.

Naturally, if no insertion or deletion can be performed but both insertion and deletion contexts have the same constructor at their root, we want to recognize this constructor as part of the spine of the alignment, and continue to align its fields pairwise.

$$\text{Spn} :: (\text{CompoundCnstr } \kappa \text{ fam } x) \Rightarrow \text{SRep } (\text{Al } \kappa \text{ fam } f) (\text{Rep } at) \rightarrow \text{Al } \kappa \text{ fam } f at$$

The Spn inside an alignment does not need to preserve metavariable scoping. Consequently, it can be pushed closer to the leaves uncovering as many copies as possible. When no Ins , Del or Spn can be used, we must resort to recording a unclassified modifi-

cation, of type $f\ at$. Most of the times f will be simply $Chg\ \kappa\ fam$, but we will be needing to add some extra information in the leaves of an alignment later. Moreover, keeping the f a parameter turns Al into a functor which enables us to map over it easily.

$Mod :: f\ at \rightarrow Al\ \kappa\ fam\ f\ at$

Take an alignment $a = Mod\ (Chg\ \#_x\ \#_x)$. Does a represent a copy or is x contracted or duplicated? Because metavariables are scoped globally within an alignment, we can only distinguish between copies and duplications by traversing the entire alignment and recording the arity of x . Yet, it is an important distinction to make. A copy synchronizes with anything whereas a contraction needs to satisfy additional constraints. Therefore, we will identify copies and permutations directly in the alignment to aid the merge function, later. Let $c = Chg\ \#_x\ \#_y$ where both x and y occur twice in their global scope: once in the deletion context and once in the insertion context. We say c is a copy when $x \equiv y$ and a permutation when $x \not\equiv y$. These are the last two constructors of Al .

$Cpy :: Metavar\ \kappa\ fam\ at \rightarrow Al\ \kappa\ fam\ f\ at$

$Prm :: Metavar\ \kappa\ fam\ at \rightarrow Metavar\ \kappa\ fam\ at \rightarrow Al\ \kappa\ fam\ f\ at$

Equipped with a definition for alignments, we move on to defining $alignChg$. Given a change c , the first step of $alignChg\ c$ is checking whether the root of c_{del} (resp. c_{ins}) can be deleted (resp. inserted). A deletion (resp. insertion) of an occurrence of a constructor X can be performed when all fields of X at this occurrence are *rigid* trees with the exception of a single recursive field – recall *rigid* trees contains no metavariables. If we can delete the root, we flag it as such and continue through the recursive *non-rigid* field. If we cannot perform a deletion at the root of c_{del} nor an insertion at the root of c_{ins} but they are constructed with the same constructor, we identify the constructor as being part of the alignments' spine. If c_{del} and c_{ins} do not have the same constructor at the root, nor are copies or permutations, we finally fallback and flag an unclassified modification.

To check whether constructors can be deleted or inserted efficiently, we must annotate rigidity information throughout our trees. The *IsRigid* datatype captures whether a tree contains any metavariables or not and is placed in every node of a tree with the *annotRigidity* function.

type *IsRigid* = *Const Bool*

annotRigidity :: Holes\ \kappa\ fam\ h\ x \rightarrow HolesAnn\ \kappa\ fam\ IsRigid\ h\ x

After annotating the trees with rigidity information, we extract the zippers that witness potential insertions or deletions. This is done by the *hasRigidZipper* function, which first extracts *all* possible zippers from the root then checks whether one of them have all of its fields marked rigid except for the focus of the zipper. If we find such a zipper, we return it wrapped in a *Just*. When a rigid zipper exists it is unique by definition, hence

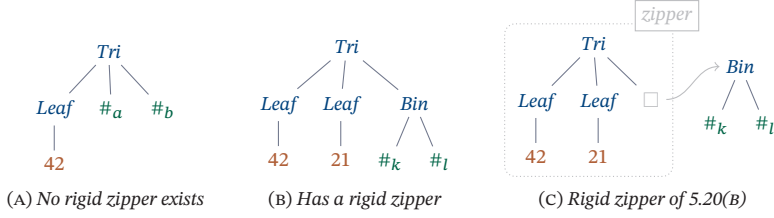


FIGURE 5.20: Examples to *hasRigidZipper* and their return values where applicable.

there is no choice involved in detecting insertions and deletions, which keeps our algorithms efficient and deterministic.

Figure 5.20 exemplifies two possible arguments to *hasRigidZipper*. The tree in Figure 5.20(A) has three possible zippers: focusing on either of its recursive positions. Neither of them, however, would have all its subtrees rigid except the focus point. Figure 5.20(B) on the other hand has one of its zippers (the one with focus on *Bin* $\#_k$ $\#_l$, Figure 5.20(C)) rigid, that is, none of the trees within the zipper has any metavariables. We omit the full implementation of *hasRigidZipper* but invite the interested reader to check *Data.HDiff.Diff.Align* in the source code (Appendix A).

Checking for deletions, then, can be easily done by first checking whether the root has a rigid zipper. If so, we can flag the deletion. In the excerpt of *alD* below, should *d* be the tree in Figure 5.20(B), *focus* would be *Bin* $\#_k$ $\#_l$, which is the single *non-rigid* recursive subtree of *d*.

```
alD d i = case hasRigidZipper d of
  Just (Zipper zd focus) → Del zd (continueAligning focus i)
```

The complete *alD* is more involved. For one, we must check whether *i* also has a rigid zipper. When both *d* and *i* have rigid zippers, we must check whether they are the same constructor and, if so, mark it as part of the spine instead. The *al* function encapsulates the *alD* above and is shown in Figure 5.21. A call to *al* will attempt to extract deletions, then insertions, then finally falling back to copies, permutations, modifications or recursively calling itself inside spines.

To compute an alignment, then, we start computing the multiset of variables used throughout a patch, annotate the deletion and insertion context with *IsRigid* and pass everything to the *al* function.

```
alignChg :: Chg κ fam at → Al κ fam (Chg κ fam) at
alignChg c@(Chg d i) = al (chgVargs c) (annotRigidity d) (annotRigidity i)
```

```

type Aligner  $\kappa$  fam = HolesAnn  $\kappa$  fam IsStiff (Metavar  $\kappa$  fam) t
    → HolesAnn  $\kappa$  fam IsStiff (Metavar  $\kappa$  fam) t
    → Al  $\kappa$  fam (Chg  $\kappa$  fam t)

al :: Map Int Arity → Aligner  $\kappa$  fam
al vars d i = alD (alS vars (al vars)) d i
where
    -- Try deleting many; try inserting one; decide whether to delete,
    -- insert or spn in case both Del and Ins are possible. Fallback to
    -- inserting many.
    alD :: Aligner  $\kappa$  fam → Aligner  $\kappa$  fam
    alD f d i = case hasRigidZipper d of -- Is the root a potential deletion?
        Nothing      → alI f d i
        -- If so, we must check whether we also have a potential insertion.
        Just (Zipper zd rd) → case hasRigitZipper i of
            Nothing      → Del (Zipper zd (alD f rd i))
            Just (Zipper zi ri) → case zipSZip zd zi of -- are zd and zi the same?
                Just res → Spn $ plug (zipperMap Mod res) (alD f rd ri)
                Nothing → Del (Zipper zd (Ins (Zipper zi (alD f rd ri))))

    -- Try inserting many; fallback to parametrized action.
    alI :: Aligner  $\kappa$  fam → Aligner  $\kappa$  fam
    alI f d i = case hasRigidZipper i of
        Nothing      → f d i
        Just (Zipper zi ri) → Ins (Zipper zi (alI f d ri))

    -- Try extracting spine and executing desired action
    -- on the leaves; fallback to deleting; inserting then modifying
    -- if no spine is possible.
    alS :: Map Int Arity → Aligner  $\kappa$  fam → Aligned  $\kappa$  fam
    alS vars f d@(Roll' _ sd) i@(Roll' _ si) =
        case zipSRep sd si of
            Nothing → alMod vars d i
            Just r  → Spn (repMap (uncurry' f) r)
    syncSpine vars _ d i = alMod vars d i

    -- Records a modification, copy or permutation.
    alMod :: Map Int Arity → Aligned  $\kappa$  fam
    alMod vars (Hole' _ vd) (Hole' _ vi) =
        -- are both vd and vi with arity 2?
        | all (≡ Just 2 ∘ flip lookup vars) [metavarGet vd, metavarGet vi]
        = if vd ≡ vi then Cpy vd else Prm vd vi
        | otherwise
        = Mod (Chg (Hole vd) (Hole vi))
    alMod _ d i = Mod (Chg d i)

```

FIGURE 5.21: Complete definition of *al*.

Forgetting the information computed by *alignChg* is trivial, as shown in the *disalign* function sketched below. It converts a *Al* back into a *Chg* in the expected way: it plugs deletion and insertion zippers and distributes the constructors in the spine into both deletion and insertion contexts and translates *Cpy* and *Prm* as expected.

```
disalign :: Al κ fam (Chg κ fam) at → Chg κ fam at
disalign (Del (Zipper del rest)) =
  let Chg d i = disalign rest in Chg (Roll (plug (cast del) d) i)
disalign ...
```

Distributing an outer spine through an alignment is trivial. All we must do is place all the constructors of the outer as constructors belonging to the alignment's spine, *Spn*.

```
alDistr :: PatchAl κ fam at → Al κ fam (Chg κ fam) at
```

Finally, computing aligned patches from locally-scoped patches is done by mapping over the outer spine and aligning the changes individually, then we make a pass over the result and issue copies for opaque values that appear on the alignment's inner spine.

```
align :: Patch κ fam at → PatchAl κ fam at
align = fst ∘ align'
```

The auxiliary function *align'* returns the successor of the last issued name to ensure we can easily produce fresh names later on, if need be. Once again, a technicality of handling names explicitly. Note that *align* introduces information, namely, new metavariables that represent copies over opaque values that appear on the alignment's spine. This means that mapping *disalign* to the result of *align* will *not* produce the same result. Alignments and changes are *not* isomorphic.

```
align' :: Patch κ fam at → (PatchAl κ fam at, Int)
align' p = flip runState maxv $ holesMapM (alRefineM cpyPrims ∘ alignChg vars) p
where vars = patchVars p
      maxv = maybe 0 id (lookupMax vars)
```

The *cpyPrims* above issues a *Cpy i*, for a fresh name *i* whenever it sees a modification with the form *Chg (Prim x) (Prim y)* with $x \equiv y$. The *alRefineM f* applies a function in the leaves of the *Al* and has type:

```
alRefineM :: (Monad m) ⇒ (∀ x . fx → m (Al κ fam g x))
→ Al κ fam f ty → m (Al κ fam g ty)
```

This process of computing alignments showcases an important aspect of our well-typed approach: the ability to access type-level information in order to compute zippers

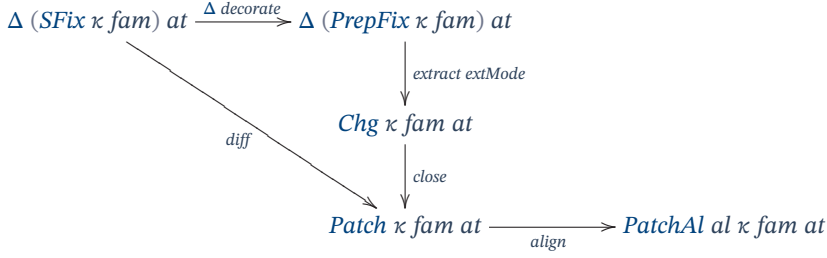


FIGURE 5.22: Conceptual pipeline of the design space for the *diff* function. $\Delta f x$ denotes $(f x, f x)$

and understand deletions and insertions of a single layer in a homogeneous fashion – the type that results from the insertion or deletion is the same type that is expected by the insertion or deletion.

5.2.4 SUMMARY

In Section 5.2 we have seen how *Chg* represents an unrestricted tree-matching, which can later be translated into isolated, well-scoped, fragments connected through an outer spine and making up a *Patch*. Finally, we have seen how to extract valuable information from well-scoped fragments about which constructors have been deleted, inserted or still belong to an inner spine, giving rise to alignments. This representation is a mix of local and global alignments. The outer spine is important to isolate a large change into smaller chunks, independent of one another.

The *diff* function produces a *Patch* instead of a *PatchAl* to keep it consistent with our previously published work [76], but also because it is easier to manage calls to *align* where they are directly necessary, since *align* produces fresh variables and this can require special attention to keep names from being shadowed.

In fact, the *diff* function could be any path in the diagram portrayed in Figure 5.22. There is no *right* choice as this depends on the specific application in question. For our particular case of pursuing a synchronization function, we require all the information up to *PatchAl*.

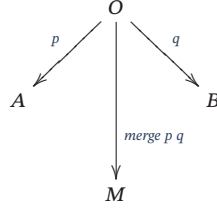


FIGURE 5.23: *Span of patches, p (transforms O into A) and q (transforms O into B). Both patches have a common element O in their domain. The patch $\text{merge } p \ q$ applies to this common ancestor O and can be thought of as the union of the changes of p and q .*

5.3 MERGING ALIGNED PATCHES

In this section we will be exploring a synchronization algorithm for aligned patches, witnessed by the *merge* function, declared below, which receives two *aligned* patches p and q that make a span – that is, have at least one common element in their domain. The result of *merge* $p \ q$ is a patch that might contain conflicts, denoted by *PatchC*, whenever both p and q modify the same subtree in two distinct ways. If p and q do *not* make a span *merge* $p \ q$ returns *Nothing*. Figure 5.23 illustrates a span of patches p and q and their merge which is supposed to be applied to their common ancestor producing a tree which combines the modifications performed by p and q , when possible.

merge :: *PatchAl* κ *fam at* \rightarrow *PatchAl* κ *fam al* \rightarrow *Maybe* (*PatchC* κ *fam at*)

Recall our patches consist of a spine which leads to locally-scoped alignments, which in turn have an inner spine that ultimately leads to changes. The distinction between the *outer* spine and the spine inside the alignments is the scope. Consequently, we can map a pure function over the outer spine without having to carry information about local scopes to the next call. When manipulating the *inner* spine, however, we must keep track of which variables have or have not been declared or used. Take the example in Figure 5.24, that merges patches p (Figure 5.24(A)) and q (Figure 5.24(B)) to produce a new patch (Figure 5.24(C)). While synchronizing the left child of each root, we discover that the tree located at (or, identified by) $\#_x$ was *Leaf* 42. We must remember this information since we will encounter $\#_x$ again and must ensure that it matches with its previously discovered value in order to perform the contraction. When we finish synchronizing the left child of the root, though, we can forget about $\#_x$ since well-scopedness of alignments guarantees $\#_x$ will not appear elsewhere.

It helps to think about metavariables in a change as a unique identifier for a subtree in the source. For example, if one change modifies a subtree x into a different subtree

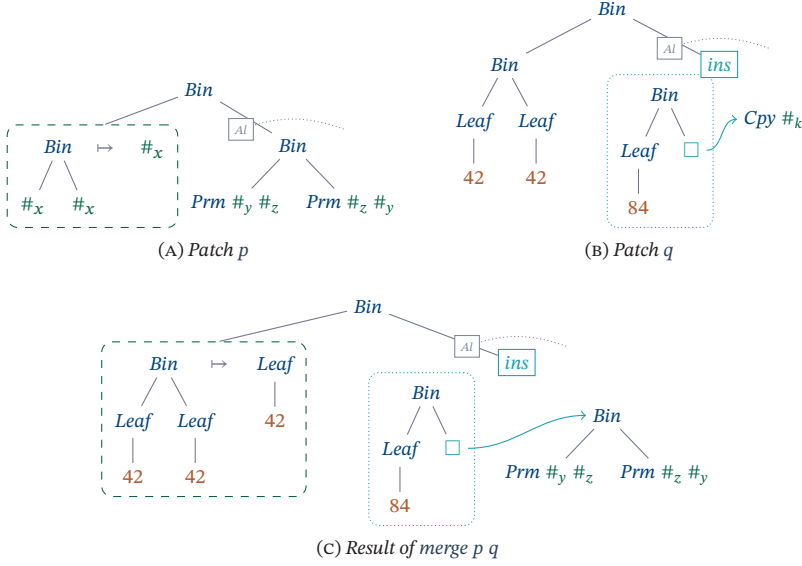


FIGURE 5.24: Example of a simple synchronization

x' , but some other change moves x , identified by $\#_x$, to a different location in the tree, the result of synchronizing these should be the transport of x' into the new location – which is exactly where $\#_x$ appears in the insertion context. The example in Figure 5.25 illustrates this very situation: the source tree identified by $\#_x$ in the deletion context of Figure 5.25(B) was changed, by Figure 5.25(A), from *Leaf* 42 into *Leaf* 84. Since p altered the content of a subtree, but q altered its location, they are *disjoint* – they alter different aspects of the common ancestor. Hence, the synchronization is possible and results in Figure 5.25(C).

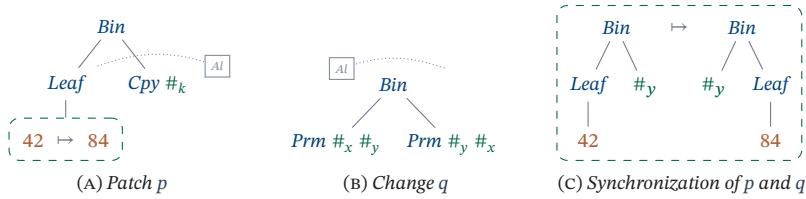


FIGURE 5.25: Example of a simple synchronization.

Given then two aligned patches, the *merge* $p\ q$ function below will map over the common prefix of the spines of p and q , captured by their least-general-generalization and produce a patch that might contain conflicts inside. In the actual implementation we receive two patches and align them inside *merge*, as this helps ensuring they will have a disjoint set of names.

```

merge :: PatchAl  $\kappa$  fam at  $\rightarrow$  PatchAl  $\kappa$  fam at  $\rightarrow$  Maybe (PatchC  $\kappa$  fam at)
merge oa ob = holesMapM (uncurry' go) (lgg oa ob)
  where go :: Holes  $\kappa$  fam (Al  $\kappa$  fam) at  $\rightarrow$  Holes  $\kappa$  fam (Al  $\kappa$  fam) at
         $\rightarrow$  Maybe (Sum (Conflict  $\kappa$  fam) (Chg  $\kappa$  fam) at)
        go ca cb = mergeAl (alDistr ca) (alDistr cb)

```

A conflict, defined below, contains a label identifying which branch of the merge algorithm issued it and the two alignments that could not be synchronized. Conflicts are issued whenever we were not able to reconcile the alignments in question. This happens either when we cannot detect that two edits to the same location are non-interfering or when two edits to the same location in fact interfere with one another. Putting it differently, conflicts might contain false positives where edits could have been automatically reconciled. The *PatchC* datatype encodes patches which might contain conflicts inside.

```

data Conflict  $\kappa$  fam at = Conflict String (Al  $\kappa$  fam at) (Al  $\kappa$  fam at)
type PatchC  $\kappa$  fam at = Holes  $\kappa$  fam (Sum (Conflict  $\kappa$  fam) (Chg  $\kappa$  fam)) at

```

Merging has a large design space. In what follows we will discuss our initial exploration and prototype algorithm, which was driven by practical experiments (Chapter 6).

The *mergeAl* function is responsible for synchronizing alignments and is where most of the work happens. In broad strokes, it is similar to synchronizing *Patch_{st}*'s (Section 4.2): insertions are preserved as long as they do not happen simultaneously. Deletions must be *applied* before continuing and copies are the identity of synchronization. In the current setting, however, we also have permutations and arbitrary changes to look at. The general conducting line of our synchronization algorithm is to first record how each subtree was modified and then instantiate these modifications in a later phase. Traversing the patches simultaneously whilst constructing substitutions would not suffice since the order which metavariables appear in each context can be drastically different. This would require us to start over every time we discovered new information on the current traversal, yielding a very slow merging algorithm.

Let us look at an example, illustrated in Figure 5.26. We start identifying we are in a situation where both *diff* $o\ a$ and *diff* $o\ b$ are spines, that is, they copy the same constructor at their root. Recursing pairwise through their children, we see a permutation versus a copy. Since a copy is the identity element, we return the permutation. On the right we see another spine versus an insertion, but since the insertion represents new information, it must be preserved. Finally, inside the insertion we see another copy,

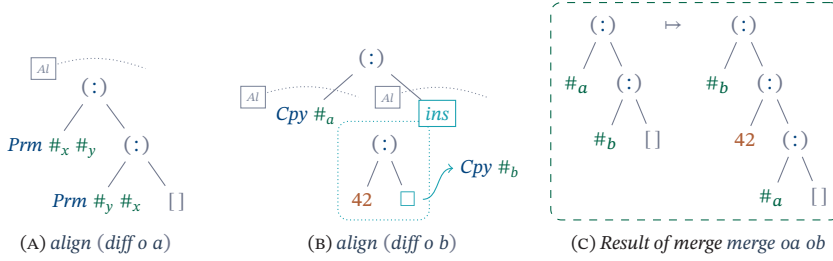


FIGURE 5.26: Example merge of two simple patches.

which means that the spine should be preserved as is. The resulting patch can be seen in Figure 5.26(c).

We keep track of the equivalences we discover in a state monad. The instantiation of metavariables will be stored under *inst* and the list of tree equivalences will be stored under *eqs*.

```
data MergeState  $\kappa$  fam = MergeState
  { inst :: Map (Exists (Metavar  $\kappa$  fam)) (Exists (Chg  $\kappa$  fam))
  , eqs  :: Map (Exists (Metavar  $\kappa$  fam)) (Exists (HolesMV  $\kappa$  fam))
  }
```

It is important to keep track of equivalences in *eqs*. Say, for example, we are to merge two changes that were left as *unclassified* by our alignment algorithm. Naturally, their deletion contexts must be unifiable, yielding a series of equivalences between their metavariables but since we do not possess information about exactly how each of those metavariables were transformed, we cannot register how they changed in *inst*. Figure 5.27 provides a simple such example. When unifying the deletion contexts of Figure 5.27(A) and Figure 5.27(B), we learn that $\{\#_x \equiv \text{Leaf } 42, \#_a \equiv \#_x, \#_b \equiv \#_y\}$, which enable us to conclude both changes are compatible and perform the same action modulo a contraction and can be merged, yielding Figure 5.27(C).

Conflicts and errors stemming from the arguments to *mergeAl* not forming a span will be distinguished by the *MergeErr* datatype, below. We also define auxiliary functions to raise each specific error in a computation inside the *Except* monad.

```
data MergeErr = NotASpan | Conf String
throwConf lbl  = throwError (Conf lbl)
throwNotASpan = throwError NotASpan
```

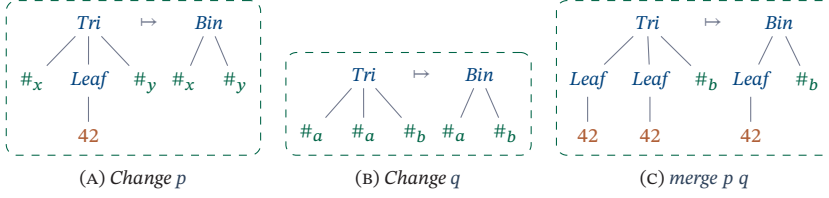



FIGURE 5.27: Merging arbitrary changes requires knowledge of equivalences between metavariables and trees.

The *mergeAl* function is defined as a wrapper around *mergeALM*, which is defined in terms of the *MergeM* monad to help carry around the necessary state and raises errors through the *Except* monad.

```

type MergeM  $\kappa$  fam = StateT (MergeState  $\kappa$  fam) (Except MergeErr)
mergeAl :: Aligned  $\kappa$  fam x  $\rightarrow$  Aligned  $\kappa$  fam x
       $\rightarrow$  Maybe (Sum (Conflict  $\kappa$  fam) (Chg  $\kappa$  fam) x)
mergeAl x y = case runExceptT (evalStateT (mergeALM p q) mrgStEmpty) of
  Left NotASpan  $\rightarrow$  Nothing
  Left (Conf err)  $\rightarrow$  Just (InL (Conflict err p q))
  Right r         $\rightarrow$  Just (InR (disalign r))

```

Finally, the *mergeALM* function maps over both alignments that we wish to merge and collects all the constraints and observations. It then attempts to split these constraints and observations into two maps: (A) a deletion map that contains information about what a subtree identified by a metavariable *was*; and (B) an insertion map that identifies what said metavariable *became*. If it is possible to produce these two idempotent substitutions, it then makes a second pass computing the final result.

```

mergeALM :: Al  $\kappa$  fam at  $\rightarrow$  Al  $\kappa$  fam at  $\rightarrow$  MergeM  $\kappa$  fam (Al  $\kappa$  fam at)
mergeALM p q = do phase1  $\leftarrow$  mergePhase1 p q
              info  $\leftarrow$  get
              case splitDelInsMap info of
                Left _  $\rightarrow$  throwConf "failed-contr"
                Right di  $\rightarrow$  alignedMapM (mergePhase2 di) phase1

```

FIRST PHASE. The first pass is computed by the *mergePhase1* function, which will populate the state with instantiations and equivalences and place values of type *Phase2* in-place in the alignment. These values instruct the second phase on how to proceed on that particular location. In between these phases of the merge algorithm we must pro-

cess the information we gathered into a deletion and an insertion map, which enables us to understand how subtrees were modified. This is done by the *splitDelInsMap* function and will be discussed in more detail later. First, we look into how the first pass instantiates metavariables and registers equivalences.

The *mergePhase1* function receives two alignments and produces a third alignment with instructions for the *second phase*. These instructions can be instantiating a change, with *P2Instantiate*, which might include a context to ensure for some consistency predicates. Or checking that two changes are α -equivalent after they have been instantiated.

```
data Phase2  $\kappa$  fam at where
  P2Instantiate :: Chg  $\kappa$  fam at  $\rightarrow$  Maybe (HolesMV  $\kappa$  fam at)  $\rightarrow$  Phase2  $\kappa$  fam at
  P2TestEq      :: Chg  $\kappa$  fam at  $\rightarrow$  Chg  $\kappa$  fam at  $\rightarrow$  Phase2  $\kappa$  fam at
```

Deciding which instruction should be performed depends on the structure of the alignments under synchronization, and is done by the *mergePhase1* function, whose cases will be discussed one by one, next.

```
mergePhase1 :: Al  $\kappa$  fam x  $\rightarrow$  Al  $\kappa$  fam x
              $\rightarrow$  MergeM  $\kappa$  fam (Al'  $\kappa$  fam (Phase2  $\kappa$  fam) x)
mergePhase1 p q = case (p, q) of
  (Cpy _, _)  $\rightarrow$  return (Mod (P2Instantiate (disalign q)))
  (_, Cpy _)  $\rightarrow$  return (Mod (P2Instantiate (disalign p)))
```

The first cases we have to handle are copies, shown above, which should be the identity of synchronization. That is, if p is a copy, all we need to do is modify the tree according to q at the current location. We might need to refine q according to other constraints we discovered in other parts of the alignment in question, so the final instruction is to *instantiate* the *Chg* that comes from forgetting the alignment q . Recall *disalign* maps alignments back into changes.

Next we look at permutations, which are almost copies in the sense that they do not modify the *content* of the tree, but they modify the *location*. We distinguish the case where both patches permute the same tree versus the case where one patch permutes the tree but the other changes its contents.

```
(Prm x y, Prm w z)  $\rightarrow$  Mod <$> mrgPrmPrm x y w z
(Prm x y, _)        $\rightarrow$  Mod <$> mrgPrm x y (disalign q)
(_, Prm x y)        $\rightarrow$  Mod <$> mrgPrm x y (disalign p)
```

When merging two permutations, *Prm* $\#_x \#_y$ against *Prm* $\#_w \#_z$, for example, we know that $\#_x$ and $\#_w$ must refer to the same subtree, hence we register their equivalence. But since the two changes permuted the same tree, we can only synchronize them if they were permuted to the *same place*. We can only assert that if we also get an equivalence

between $\#_y$ and $\#_z$. Consequently, we issue a *P2TestEq* and establish whether the two permutations can be merged or not at the end of process.

```

mrgPrmPrm :: Metavar  $\kappa$  fam  $x \rightarrow$  Metavar  $\kappa$  fam  $x$ 
            $\rightarrow$  Metavar  $\kappa$  fam  $x \rightarrow$  Metavar  $\kappa$  fam  $x$ 
            $\rightarrow$  MergeM  $\kappa$  fam (Phase2  $\kappa$  fam  $x$ )
mrgPrmPrm  $x\ y\ w\ z =$  onEqvs ( $\lambda eqs \rightarrow$  substInsert eqs  $x$  (Hole  $w$ ))
            $\gg$  return (P2TestEq (Chg (Hole  $x$ ) (Hole  $y$ )) (Chg (Hole  $w$ ) (Hole  $z$ )))

```

If we are merging one permutation with something other than a permutation, however, we know one change modified the location of a tree, whereas another potentially modified its contents. All we must do is record that the tree identified by $\#_x$ was modified according to c . After we have made one entire pass over the alignments being merged, we must instantiate the permutation with the information we discovered – the $\#_x$ occurrence in the deletion context of the permutation will be c_{del} , potentially simplified or refined. The $\#_y$ appearing in the insertion context of the permutation will be instantiated with whatever we come to discover about it later. We know there *must* be a single occurrence of $\#_y$ in a deletion context because the alignment flagged it as a permutation.

```

mrgPrm :: Metavar  $\kappa$  fam  $x \rightarrow$  Metavar  $\kappa$  fam  $x \rightarrow$  Chg  $\kappa$  fam  $x$ 
         $\rightarrow$  MergeM  $\kappa$  fam (Phase2  $\kappa$  fam  $x$ )
mrgPrm  $x\ y\ c =$  addToInst "prm-chg"  $x\ c$ 
            $\gg$  return (P2Instantiate (Chg (Hole  $x$ ) (Hole  $y$ )) Nothing)

```

The *addToInst* function inserts the (x, c) entry in *inst* if x is not yet a member. It raises a conflict, in general, if x is already in *inst* with a value that is different from c ¹. In the call to *addToInst* in *mrgPrm*, above, it never raises a “*prm-chg*” conflict. This is because $\#_x$ and $\#_y$ are classified as a permutation – each variable occurs exactly once in the deletion and once in the insertion contexts. Therefore, it is impossible that x was already a member of *inst*.

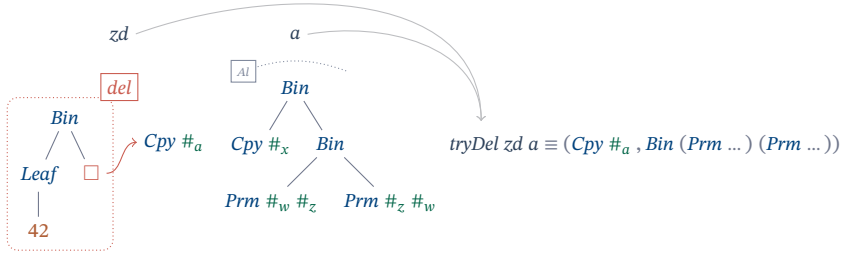
Next, we look at insertions. Interstinos must be preserved as long as they do not attempt to insert different information in the same location, otherwise we would not be able to decide which insertion should happen first.

```

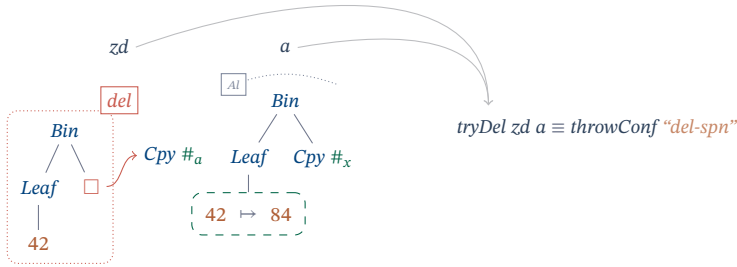
(Ins (Zipper  $z\ p'$ ), Ins (Zipper  $z'\ q'$ ))
|  $z \equiv z'$             $\rightarrow$  Ins  $\circ$  Zipper  $z <\$>$  mergePhase1  $p'\ q'$ 
| otherwise           $\rightarrow$  throwConf "ins-ins"
(Ins (Zipper  $z\ p'$ ),  $-$ )  $\rightarrow$  Ins  $\circ$  Zipper  $z <\$>$  mrgPhase1  $p'\ q$ 
( $-$ , Ins (Zipper  $z\ q'$ ))  $\rightarrow$  Ins  $\circ$  Zipper  $z <\$>$  mrgPhase1  $p\ q'$ 

```

¹ Instead of forbidding values different than c , we could think to check whether the two different values can be merged. This would incur other difficulties and is left as future work.



(A) Call to *tryDel* succeeds. The *Bin* at the root can be deleted as it only overlaps with copies. *tryDel* returns the focus of the deletion and the part of the alignment *a* that overlaps with it.



(B) Call to *tryDel* fails. Although the *Bin* at the root could be deleted, the alignment *a* is changing the 42 present in the leaf. This is a conflict.

FIGURE 5.28: Two example calls to *tryDel*.

Deletions must be preserved and *executed*. That is, if one patch deletes a constructor but the other modifies the fields of the constructor, we must first ensure that none of the deleted fields have been modified but the deletion should be preserved in the merge. The *tryDel* function attempts to execute the deletion of a zipper over an alignment, and, if successful, returns the pair of alignments we should continue to merge. It essentially overlaps the deletion zipper with *a* and observes whether *a* performs no modifications anywhere except on the focus of the zipper. When its not possible to execute the deletion we can continue. Figure 5.28 illustrates some example calls to *tryDel*, whose complete generic definition is shown in Figure 5.29.

$$\begin{aligned}
 (\text{Del } zp @ (\text{Zipper } z \text{ } _), \text{ } _) &\rightarrow \text{Del} \circ \text{Zipper } z \text{ } \langle \$ \rangle (\text{tryDel } zp \text{ } q \gg \text{uncurry mrgPhase1}) \\
 (\text{ } _, \text{Del } zq @ (\text{Zipper } z \text{ } _)) &\rightarrow \text{Del} \circ \text{Zipper } z \text{ } \langle \$ \rangle (\text{tryDel } zq \text{ } p \gg \text{uncurry mrgPhase1})
 \end{aligned}$$

Note that since *merge* is supposed to be symmetric, we can freely swap the order of arguments. Although we never got to proving this formally, our QuickCheck tests were encouraging that this is the case.

```

tryDel :: Zipper (CompoundCnstr κ fam x) (SFix κ fam) (Al κ fam (Chg κ fam)) x
  → Al κ fam (Chg κ fam) x
  → MergeM κ fam (Al κ fam (Chg κ fam) x, Al κ fam (Chg κ fam) x)
tryDel (Zipper z h) (Del (Zipper z' h'))
  | z ≡ z'    = return (h, h')
  | otherwise = throwConf "del-del"
tryDel (Zipper z h) (Spn rep) = case zipperRepZip z rep of
  Nothing → throwNotASpan
  Just r   → case partition (exElim isInR1) (repLeavesList r) of
    ([Exists (InL Refl :*: x)], xs)
      | all isCpyL1 xs → return (h, x)
      | otherwise     → throwConf "del-spn"
    _                  → error "unreachable; zipRepZip invariant"
tryDel (Zipper _ _) = throwConf "del-mod"

```

FIGURE 5.29: Complete generic definition of the *tryDel* function.

Next we have spines versus modifications. Intuitively, we want to match the deletion context of the change against the spine and, when successful, return the result of instantiating the insertion context of the change.

$$\begin{aligned}
 (\text{Mod } p', \text{Spn } q') &\rightarrow \text{Mod} \langle \$ \rangle \text{ mrgChgSpn } p' \ q' \\
 (\text{Spn } p', \text{Mod } q') &\rightarrow \text{Mod} \langle \$ \rangle \text{ mrgChgSpn } q' \ p'
 \end{aligned}$$

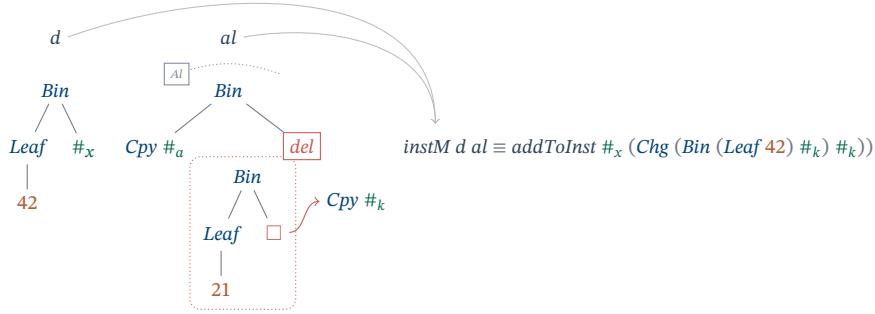
The *mrgChgSpn* function, below, matches the deletion context of the *Chg* against the spine and returns a *P2Instantiate* instruction. The instantiation function *instM*, illustrated in Figure 5.30 and defined in Figure 5.31, receives a deletion context and an alignment and attempts to assign the variables in the deletion context to changes inside the alignment. This is only possible, though, when the modifications in the spine occur *further* from the root than the variables in the deletion context. Otherwise, we have a conflict where some constructors flagged for deletion are also marked as modifications.

```

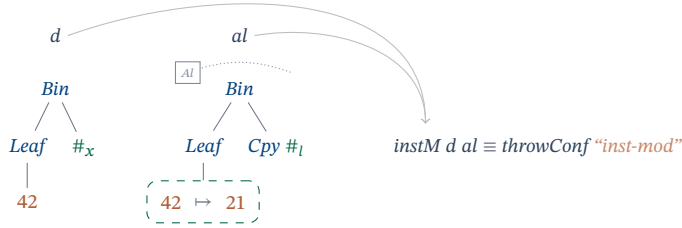
mrgChgSpn :: (CompoundCnstr κ fam x) ⇒ Chg κ fam x → SRep (Al κ fam) (Rep x)
  → MergeM κ fam (Phase2 κ fam x)
mrgChgSpn p@(Chg dp _) spn = do
  instM dp (Spn spn)
  return (P2Instantiate p (Just (disalign (Spn spn))ins))

```

The *Just* in the return value above indicates that we must check that we will not introduce extra duplications. Figure 5.32 illustrates a case where failing to perform this check would result in an erroneous duplication of the value 2. Matching the deletion context of *chg* = *Chg* #_c (#_a : #_c) against the spine *spn* = *Spn* (Cpy #_o : Chg #_z (#_x : #_z))



(A) Call to `instM` succeeds and registers that the subtree identified by `#x` has had its left child deleted, according to the alignment.



(B) Call to `instM` returns a conflict. The deletion context, `d`, wants to match against the value `42` but the alignment modifies it.

FIGURE 5.30: Two example calls to `instM`.

```

instM :: (All Eq κ) ⇒ HolesMV κ fam at → Al κ fam at → MergeM κ fam ()
instM _ (Cpy _) = return ()
instM (Hole v) a = addToInst "inst-contr" v (disalign a)
instM _ (Mod _) = throwConf "inst-mod"
instM _ (Prm _ _) = throwConf "inst-perm"
-- Del ctx and spine must form a span; cannot reference different constructors or primitives.
instM x@(Prim _) d = when (x ≠ (disalign d)del) throwNotASpan
instM (Roll r) (Spn s) = case zipSRep r s of
  Nothing → throwNotASpan
  Just res → void (repMapM (λx → uncurry' instM x >> return x) res)
instM (Roll _) (Ins _) = throwConf "inst-ins"
instM (Roll _) (Del _) = throwConf "inst-del"

```

FIGURE 5.31: Implementation of `instM`, which receives a deletion context and an alignment and attempts to instantiate the variables in the deletion context with changes in the alignment.

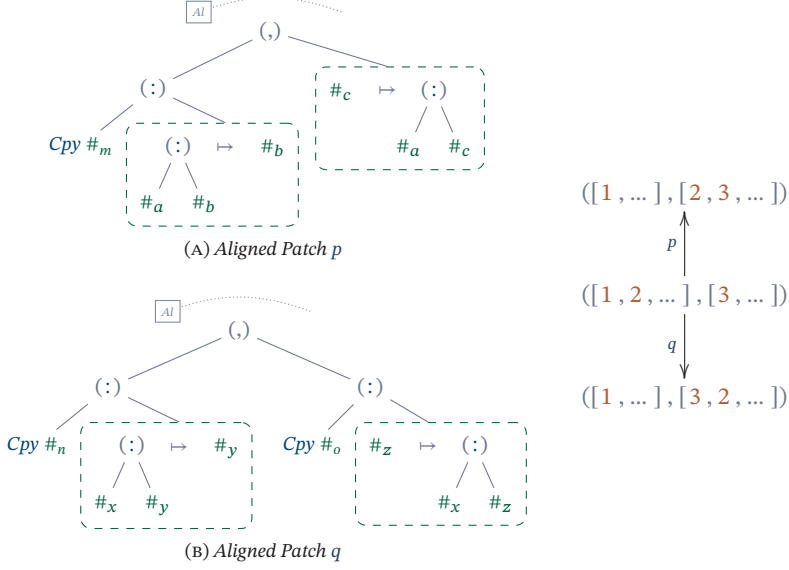
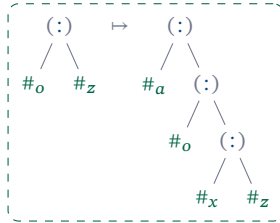


FIGURE 5.32: Example of two conflicting patches that move the same subtree into two different locations. The patches here are operating over pairs of lists.

yields $\#_c$ equal to spn , which correctly identifies that the subtree at $\#_c$ was modified according to spn . The observation, however, is that the insertion context of chg mentions $\#_a$, which is a subtree that comes from outside the deletion context of chg . If we do not perform any further check and proceed naively, we would end up substituting $\#_c$ for $(disalign\ spn)_{del}$ and for $(disalign\ spn)_{ins}$ in chg_{del} and chg_{ins} , respectively, which would result in:



Since we know $\#_x \equiv \#_a$, which was registered when merging the left hand side of $(,)$, in Figures 5.32(A) and 5.32(B), it becomes clear that $\#_a$ was erroneously duplicated. Our implementation will reject this by checking that the set of subtrees that appear in the result of instantiating chg is disjoint from the set of subtrees moved by spn .

Merging two spines is simple. We know they must reference the same constructor since the arguments to *merge* form a span. All that we have to do is recurse on the paired fields of the spines, point-wise:

```
(Spn p' , Spn q') → case zipSRep p' q' of
  Nothing → throwNotASpan
  Just r → Spn <$> repMapM (uncurry' mrgPhase1) r
```

Lastly, when the alignments in question are arbitrary modifications, we must try our best to reconcile these. We handle duplications differently from arbitrary modifications, which are easier to handle.

```
(Mod p' , Mod q') → Mod <$> mrgChgChg p' q'
```

A duplication or contraction is of the form *Chg* #_x #_y, where #_x or #_y occurs at least three times in the alignment at question. Three occurrences might seem arbitrary, but a metavariable must occur at least twice, and, when it occurs only twice the alignment algorithm would have marked it as a copy or a permutation. Merging duplications is straightforward. When either one of *p'* or *q'* above are a duplication but the other is a change, we record how the tree was changed and move on.

```
mrgChgDup :: Chg κ fam x → Chg κ fam x → MergeM κ fam (Phase2 κ fam x)
mrgChgDup dup@(Chg (Hole v) _) q' = do
  addToInst "chg-dup" v q'
  return (P2Instantiate dup Nothing)
```

Finally, if *p* and *q* are not duplications, nor any of the cases previously discussed, then the best we can do is register equivalence of their domains – recall both patches being merged must form a span – and synchronize successfully when both changes are equal.

```
mrgChgChg :: Chg κ fam x → Chg κ fam x → MergeM κ fam (Phase2 κ fam x)
mrgChgChg p' q' | isDup p' = mrgChgDup p' q'
                  | isDup q' = mrgChgDup q' p'
                  | otherwise = case unify p' del q' del of
                    Left _ → throwNotASpan
                    Right r → onEqvs (M.∪ r) >> return (P2TestEq p' q')
```

Once the first pass is done, we have collected information about how each subtree has been changed and potential subtree equivalences we might have discovered. The next step is to synthesize this information into two maps: a deletion map that informs us what a subtree *was* and an insertion map that informs us what a subtree *became*, so we can perform the *P2Instante* and *P2TestEq* instructions.

SECOND PHASE. The second phase starts with splitting *inst* and *eqvs*, which requires some attention. For example, imagine there exists an entry in *inst* that assigns $\#_x$ to *Chg* (*Hole* $\#_y$) (*42* : (*Hole* $\#_y$)), this tells us that the tree identified by $\#_x$ is the same as the tree identified by $\#_y$, and it became *42* : $\#_y$. Now suppose that $\#_x$ was duplicated somewhere else, and we come across an equivalence that says $\#_y \equiv \#_x$. We cannot simply insert this equivalence into *inst* because the merge algorithm made the decision to remove all occurrences of $\#_x$, not of $\#_y$, even though they identify the same subtree. This is important to ensure we produce patches that can be applied.

The *splitDelInsMaps* function is responsible for synthesizing the information gathered in the first pass of the synchronization algorithm. First we split *inst* into the deletion and insertion components of each of its points. Next, we partition the equivalences into rigid equivalences, of the form $(\#_v, t)$ where *t* has no holes, and non-rigid equivalences. The rigid equivalences are added to both deletion and insertion maps, but the non-rigid ones, $(\#_v, t)$, are only added when there is no information about the $\#_v$ in the map and, if $t \equiv \#_u$, we also check that there is no information about $\#_u$ in the map. Lastly, after these have been added to the map, we call *minimize* to produce an idempotent substitution we can use for phase two. If an occurs-check error is raised, this is forwarded as a conflict.

```

type Subst2  $\kappa$  fam = (Subst  $\kappa$  fam (Metavar  $\kappa$  fam) , Subst  $\kappa$  fam (Metavar  $\kappa$  fam))
splitDelInsMaps :: MergeState  $\kappa$  fam  $\rightarrow$  Either [Exists (Metavar  $\kappa$  fam)] (Subst2  $\kappa$  fam)
splitDelInsMaps (MergeState iot eqvs) = do
  let e' = splitEqvs eqvs
  d  $\leftarrow$  addEqvsAndSimpl (map (exMap  $\cdot_{\text{del}}$ ) inst) e'
  i  $\leftarrow$  addEqvsAndSimpl (map (exMap  $\cdot_{\text{ins}}$ ) inst) e'
  return (d , i)

```

After computing the insertion and deletion maps, which inform us how each identified subtree was modified, we start a second pass over the result of the first pass and execute the necessary instructions.

```

phase2 :: Subst2  $\kappa$  fam  $\rightarrow$  Phase2  $\kappa$  fam at  $\rightarrow$  MergeM  $\kappa$  fam (Chg  $\kappa$  fam at)
phase2 di (P2TestEq ca cb)      = isEqChg di ca cb
phase2 di (P2Instantiate chg Nothing) = return (refineChg di chg)
phase2 di (P2Instantiate chg (Just i)) = do
  es  $\leftarrow$  gets eqs
  case getCommonVars (substApply es chgins) (substApply es i) of
    []  $\rightarrow$  return (refineChg di chg)
    xs  $\rightarrow$  throwConf ("mov-mov " ++ show xs)

```

The *getCommonVars* computes the intersection of the variables in two *Holes*, which is used to forbid subtrees to be moved in two different ways.

Refining changes according to the inferred information is straightforward, all we must do is apply the deletion map to the deletion context and the insertion map to the insertion context.

$$\begin{aligned} \text{refineChg} &:: \text{Subst2 } \kappa \text{ fam} \rightarrow \text{Chg } \kappa \text{ fam at} \rightarrow \text{Chg } \kappa \text{ fam at} \\ \text{refineChg } (d, i) (\text{Chg } \text{del ins}) &= \text{Chg } (\text{substApply } d \text{ del}) (\text{substApply } i \text{ ins}) \end{aligned}$$

When deciding whether two changes are equal, its also important to refine them first, since they might be α -equivalent.

$$\begin{aligned} \text{isEqChg} &:: \text{Subst2 } \kappa \text{ fam} \rightarrow \text{Chg } \kappa \text{ fam at} \rightarrow \text{Chg } \kappa \text{ fam at} \rightarrow \text{Maybe } (\text{Chg } \kappa \text{ fam at}) \\ \text{isEqChg } di \text{ ca } cb &= \text{let } ca' = \text{refineChg } di \text{ ca} \\ &\quad cb' = \text{refineChg } di \text{ cb} \\ &\quad \text{in if } ca' \equiv cb' \text{ then Just } ca' \text{ else Nothing} \end{aligned}$$

The merging algorithm presented in this section is involved. It must deal with a number of corner cases and use advanced techniques to do so generically. Most of the difficulties come from having to deal with arbitrary duplications and contractions. If we instead chose to use only linear patches, that is, patches where each metavariable must be declared and used exactly once, the merge algorithm could be simplified.

5.4 DISCUSSION AND FURTHER WORK

With `hdiff` we have seen that a complete detachment from edit-scripts enables us to define a computationally efficient differencing algorithm and how the notion of *change* coupled with a simple notion of composition gives a sensible algebraic structure. The patch datatype in `hdiff` is more expressive than edit-script based approaches, as it enables us to write transformations involving arbitrary permutations and duplications. As a consequence, we have a more involved merge algorithm. For one, we cannot easily generalize our three-way merge to n -way merge. More importantly, though, there are subtleties in the algorithm that arose purely from practical necessities. Our posterior empirical evaluation (Chapter 6) does indicate that the best success ratio comes from merging linear patches – where metavariables occur exactly twice, obtained with the *Patience* extraction mode. This does suggest that the soft-spot in the design space might well be allowing arbitrary permutations, enabling a fast differencing algorithm, but forbidding arbitrary duplications and contractions, which could enable a simpler merging algorithm. Besides the merging algorithm, we will discuss a number of other important aspects that were left as future work and would need to be addressed to bring `hdiff` from a prototype to a production tool.

REFINING MATCHING AND SHARING CONTROL

The matching engine underlying `hdiff` uses hashes indiscriminately, all information under a subtree being used to compute its hash, which can be undesirable. Imagine a parser that annotates its resulting AST with source-location tokens. This means that we would not be able to recognize permutations of statements, for example, since both occurrences would have different source-location tokens and, consequently, different hashes.

This issue goes hand in hand with deciding which parts of the tree can be shared and up until which point. For example, we probably never want to share local statements outside their scope. Recall we avoided this issue by restricting whether a subtree could be shared or not based on its height. This was a pragmatic design choice that enabled us to make progress but it is a work-around at best.

Salting the hash function of *preprocess* is not an option for working around the issue of sharing control. If the information driving the salt function changes, none of the subtrees under there can be shared again. To illustrate this, suppose we push scope names into a stack with a function *intrScope* :: *SFix* κ *fam* *at* \rightarrow *Maybe String*, which would be supplied by the user. It returns a *Just* whenever the datatype in question introduces a scope. The *const Nothing* function works as a default value, meaning that the mutually recursive family in question has no scope-dependent naming. A more interesting *intrScope*, for some imaginary mutually recursive family, is given below.

```
intrScope m@(Module ...)    = Just (moduleName m)
intrScope f@(FunctionDecl ...) = Just (functionName f)
intrScope _                 = Nothing
```

With *intrScope* as above, we could instruct the *preprocess* to push module names and function names every time it traverses through one such element of the family. For example, preprocessing the pseudo-code below would mean that the hash for a `inside fib` would be computed with [“*m*”, “*fib*”] as a salt; but a `inside fat` would be computed with [“*m*”, “*fat*”] as a salt, yielding a different hash.

```
module m
  fib n = let a = 0; b = 1; ...
  fat n = let a = 0; ...
```

This will work out well for many cases, but as soon as a change altered any information that was being used as a salt, nothing could be shared anymore. For example, if we rename `module m` to `module x`, the source and destination would contain no common hashes, since we would have used [“*m*”] to salt the hashes for the source tree, but [“*x*”] for the destination, yielding different hashes.

This problem is twofold, however. Besides identifying the algorithmic means to ensure `hdiff` could be scope-aware and respect said scopes, we must also engineer an interface to enable the user to easily define said scopes. We could think of design that made use of a custom version of the `generics-simplistic` library, with an added alias for the identity functor that could receive special treatment, for example:

```
newtype Scoped f = Scoped { unScoped :: f }
data Decl = ImportDecl ...
          | FunDecl String [ParmDecl] (Scoped Body)
...

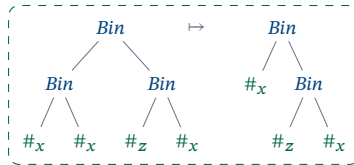
```

This would mean that when inspecting and pattern matching on *SRep* throughout our algorithms, we could treat *scoped* types differently.

We reiterate that if there is a solution to this problem, it certainly will not use a modification of the matching mechanism: if we use scope names, renamings will cause problems; if we use the order in which scopes have been seen (De Bruijn-like), permutations will cause problems. Controlling on the height of the trees and minimizing this issue was the best option to move forward in an early stage. Unfortunately, there was no time to explore how scope graphs [81] could help us here, but it is certainly a good place to start looking. It might be possible to use scope graphs to write a more intricate *close* function, that will properly break sharing where necessary, for example.

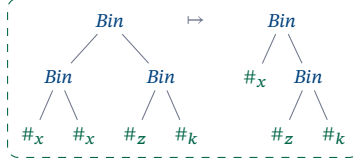
EXTRACTION METHODS, *BEST* PATCH AND EDIT-SCRIPTS

We have presented three extraction methods, which we called *NoNested*, *ProperShare* and *Patience*. Computing the diff between two trees using different extraction methods can produce different patches. Certainly there can be more extraction methods. One such example that would be interesting to implement is a refinement of *ProperShare*, aimed at breaking the sharing introduced by it. The idea was to list the metavariables that appear in the deletion and insertion context and compute the LCS between these lists. The location of copies enable us to break sharing and introduce new metavariables. For example, take the change below.



The list of metavariables in the deletion context is $[\#_x, \#_x, \#_z, \#_x]$, but in the insertion context we have $[\#_x, \#_z, \#_x]$. Computing the longest common subsequence between

these lists yields $[Del\ x, Cpy, Cpy, Cpy]$. The first *Del* suggests a contraction is really necessary, but the last copy shows that we could *break* the sharing by renaming $\#_x$ to $\#_k$, for example. This would essentially transform the change above into:



The point is that the copying of $\#_z$ can act as a synchronization point to introduce more variables, forget some sharing constraints, and ultimately enlarge the domain of our patches.

Forgetting about sharing is just one example of a different context extraction mechanism and, without a formal notion about when a patch is *better* than another, it is impossible to make a decision about which context extraction should be used. Our experimental results suggest that *Patience* yields patches that merge successfully more often, but this is far from providing a metric on patches, like the usual notion of cost does for edit-scripts.

Another interesting aspect that could have been looked at is the relation between our *Patch* datatype and traditional edit-scripts. The idea of breaking sharing above can be used to translate our patches to an edit-script. Some early experiments did show that we could use this method to compute approximations of the least-cost edit-script in linear time. Given that the minimum cost edit-script takes nearly quadratic time [11], it might be worth looking into how good an approximation we might be able to compute in linear time.

FORMALIZATIONS AND GENERALIZATIONS

Formalizing and proving properties about our *diff* and *merge* functions was also a priority. As it turns out, the extensional nature of *Patch* makes for a difficult Agda formalization, which is the reason this was left as future work.

The value of a formalization goes beyond enabling us to prove important properties. It also provides a laboratory for generalizing aspects of the algorithms. Two of those immediately jump to mind: generalizing the merge function to merge n patches and generalizing alignments insertions and deletions zippers to be of arbitrary depth, instead of a single layer. Finally, a formalization also provides important value in better understanding the merge algorithm.



EXPERIMENTS

Throughout this thesis we have presented two approaches to structural differencing. In Chapter 4 we saw `stdiff`, which although unpractical, provided us with important insights into the representation of patches. These insights and experience led us to develop `hdiff`, in Chapter 5, which improved upon the previous approach with a more efficient *diff* function at the expense of the simplicity of the merge algorithm: the *merge* function from `hdiff` is much more involved than that of `stdiff`.

In this chapter we evaluate our algorithms on real-world conflicts extracted from GitHub and analyze the results. We are interested in performance measurements and synchronization success rates, which are central factors to the applicability of structural differencing in the context of software version control.

To conduct the aforementioned evaluation we have extracted a total of 12 552 usable datapoints from GitHub. They have been obtained from large public repositories storing code written in Java, JavaScript, Python, Lua and Clojure. The choice of programming languages was motivated by the availability of parsers, with the exception of Clojure, where we borrowed a parser from a MSc thesis [36]. More detailed information about the data collection is given in Section 6.1.

The evaluation of `stdiff` has fewer datapoints than `hdiff` for the sole reason that `stdiff` requires the `generics-mrsop` library, which triggers a memory leak in GHC¹ when used with larger abstract syntax trees. Consequently, we could only evaluate `stdiff` over the Clojure and Lua subset of our dataset.

¹<https://gitlab.haskell.org/ghc/ghc/issues/17223> and <https://gitlab.haskell.org/ghc/ghc/issues/14987>

6.1 DATA COLLECTION

Collecting files from GitHub can be done with the help of some bash scripting. The overall idea is to extract the merge conflicts from a given repository by listing all commits c with more than two parents, recreating the repository at the state immediately previous to c and attempting to call `git merge` at that state.

Our script improves upon the script written by Garuffi [36] by making sure to collect the file that a human committed as the resolution of the conflict, denoted `M.lang`. To collect conflicts from a repository, then, all we have to do is run the following commands at its root.

- List each commit c with at least two parents with `git rev-list --merges`.
- For each commit c as above, let its parents be p_0 and ps ; checkout the repository at p_0 and attempt to `git merge --no-commit ps`. The `--no-commit` switch is important since it gives us a chance to inspect the result of the merge.
- Next we parse the output of `git ls-files --unmerged`, which provides us with the three *object-ids* for each file that could not be automatically merged: one identifier for the common ancestor and one identifier for each of the two diverging replicas.
- Then we use `git cat-file` to get the files corresponding to each of the *object-ids* gathered on the previous step. This yields three files, `O.lang`, `A.lang` and `B.lang`. Lastly, we use `git show` to save the file `M.lang` that was committed by a human resolving the conflict.

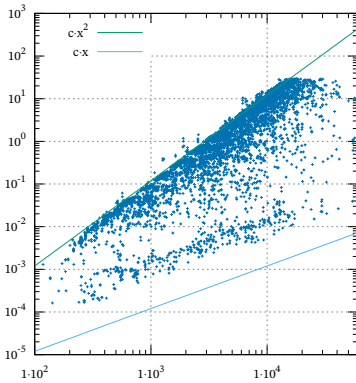
After running the steps above for a number of repositories, we end up with a list of folders containing a merge conflict that was solved manually. Each of these folders contain a span $A \leftarrow O \rightarrow B$ and a file M which is the human-produced result of synchronizing A and B . We refer the reader to the full code for more details (Appendix A). Overall, we acquired 12 552 usable conflicts – that is, we were able to parse the four files with the parsers available to us – and 2 771 conflicts where at least one file yielded a parse error. Table 6.1 provides the distribution of datapoints per programming language and displays the number of conflicts that yielded a parse error. These parse errors are an inevitable consequence of using off-the-shelf parsers on an existing dataset. The parseable conflicts have been compiled into a publicly available dataset [72].

6.2 PERFORMANCE

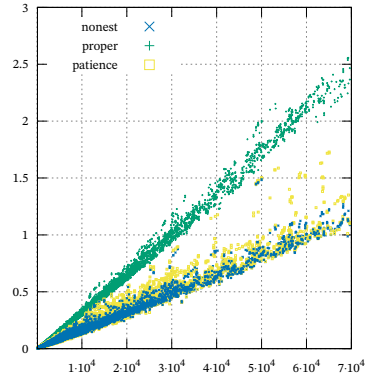
To measure the performance of the *diff* functions in both approaches we computed four patches per datapoint, namely: `diff O A`, `diff O B`, `diff O M` and `diff A B`.

Language	Repositories	Parseable Conflicts	Non-parseable Conflicts
Clojure	31	1 213	16
Java	19	2 901	851
JavaScript	28	3 392	965
Lua	27	748	91
Python	27	4 298	848
<i>Totals</i>	132	12 552	2 771

TABLE 6.1: Distribution of datapoints within our dataset [72]. The repositories were chosen manually by searching each respective language in GitHub. Our criteria for selecting repositories to mine was based on number of forks and commits, in an attempt to maximize pull requests.



(A) Runtimes from `stdiff` shown in a log-log plot. The lines illustrate the behavior of `stdiff` being between linear and quadratic



(B) Runtimes from `hdiff` shown in a linear plot.

FIGURE 6.1: Performance measurements of `stdiff` and `hdiff` differencing functions. The vertical axis represents seconds and the horizontal axis has the sum of the number of constructors in the source and destination trees.

Whilst computing patches we limited the memory usage to 8GB and runtime to 30s. If a call to *diff* used more than the available temporal and spatial resources it was automatically killed. We ran both *stdiff* and *hdiff* on the same machine, yet, we stress that the absolute values are of little interest. The real take away from this experiment is the empirical validation of the complexity class of each algorithm. The results are shown in Figure 6.1 and plot the measured runtime against the sum of the number of constructors in the source and destination trees.

Figure 6.1(A) illustrates the measured performance of the differencing algorithm in *stdiff*, our first structural differencing tool, discussed in Section 4.3.2. For *fa* and *fb* be the files being differenced, we have only timed the call to *diff fa fb* – which excludes parsing. Note that most of the time, *stdiff* exhibits a runtime proportional to the square of the input size. That was expected since it relies on a quadratic algorithm to annotate the trees and then translate the annotated trees into *Patch_{ST}* over a single pass. Out of the 8428 datapoints where we attempted to time *stdiff* in order to produce Figure 6.1(A), 913 took longer than thirty seconds and 929 used more than 8GB of memory. The rest are plotted in Figure 6.1(A). The high memory usage for particularly large examples is unsurprising. Computing a *stdiff* patch requires us to maintain and manipulate a number of singleton types and constraints.

Figure 6.1(B) illustrates the measured performance of the differencing algorithm underlying *hdiff*, discussed in Section 5.1.4. We have plotted each of the context extraction techniques described in 5.1.4.2. The linear behavior is evident and in general, an order of magnitude better than *stdiff*. We do see, however, that the proper context extraction is slightly slower than *nonest* or *patience*. Finally, only 14 calls timed-out and none used more than 8GB of memory.

Measuring performance of pure Haskell code is subtle due to its lazy evaluation semantics. We have used the *time* auxiliary function below. We based ourselves on the *timeit* package, but adapted it to fully force the evaluation of the result of the action, with the *deepseq* method and force its execution with the bang pattern in *res*, ensuring the thunk is fully evaluated.

```
time :: (NFData a) => IO a -> IO (Double, a)
time act = do t1 ← getCPUTime
             result ← act
             let ! res = result `deepseq` result
             t2 ← getCPUTime
             return (fromIntegral (t2 - t1) * 1e-12, res)
```

Language	<i>success</i>	(ratio)	<i>mdif</i>	(ratio)	total ratio	<i>conf</i>	<i>t/o</i>
Clojure	184	(0.15)	211	(0.17)	0.32	818	0
Java	978	(0.34)	479	(0.16)	0.5	1 443	1
JavaScript	1 046	(0.30)	274	(0.08)	0.38	2 062	10
Lua	185	(0.25)	101	(0.14)	0.39	462	0
Python	907	(0.21)	561	(0.13)	0.34	2 829	1
<i>Total</i>	3 300	(0.26)	1626	(0.13)	0.39	7 614	12

TABLE 6.2: *Best synchronization success rate per language. No apply-fail was encountered in the entire dataset and the number of timeouts was negligible.*

6.3 SYNCHRONIZATION

While the performance measurements provide some empirical evidence that `hdiff` is indeed linear, the synchronization experiment, discussed in this section, aims at establishing a lower bound on the number of conflicts that could be solved in practice.

The synchronization experiment consists of attempting to merge the $A \leftarrow O \rightarrow B$ span for every datapoint. If `hdiff` produces a patch with no conflicts, we apply it to O and compare the result against M , which was produced by a human. There are four possible outcomes, three of which we expect to see and one that would indicate a more substantial problem. The three outcomes we expect to see are: *success*, which indicates the merge was successful and was equal to that produced by a human; *mdif* which indicates that the merge was successful but different from the manual merge; and finally *conf* which means that the merge was unsuccessful. The other possible outcome comes from producing a patch that *cannot* be applied to O , which is referred to as *apply-fail*. Naturally, timeout or out-of-memory exceptions can still occur and fall under *other*. The merge experiment was capped at 45 seconds of runtime and 8GB of virtual memory.

The distinction between *success* and *mdif* is important. Being able to merge a conflict but obtaining a different result from what was committed by a human does not necessarily imply that either result is wrong. Developers can perform *more or fewer* modifications when committing M . For example, Figure 6.2 illustrates an example distilled from our dataset which the human performed an extra operation when merging, namely adapting the *sheet* field of one replica. It can also be the case that the developer made a mistake which was fixed in a later commit. Therefore, a result of *mdif* in a datapoint does not immediately indicate the wrong behavior of our merging algorithm. The success rate, however, provides us with a reasonable lower bound on the number of conflicts that can be solved automatically, in practice.

Language	Mode	Height	success	(ratio)	mdif	(ratio)	conf	t/o
Clojure	<i>Patience</i>	1	184	(0.15)	211	(0.17)	818	0
	<i>NoNested</i>	3	149	(0.12)	190	(0.16)	874	0
	<i>ProperShare</i>	9	92	(0.08)	84	(0.07)	1 037	0
Java	<i>Patience</i>	1	978	(0.34)	479	(0.16)	1 443	1
	<i>NoNested</i>	3	924	(0.32)	509	(0.18)	1 467	1
	<i>ProperShare</i>	9	548	(0.19)	197	(0.07)	2 155	1
JavaScript	<i>Patience</i>	1	1 046	(0.30)	274	(0.08)	2 062	10
	<i>NoNested</i>	3	991	(0.29)	273	(0.08)	2 124	4
	<i>ProperShare</i>	9	748	(0.22)	116	(0.03)	2 508	20
Lua	<i>Patience</i>	3	185	(0.25)	101	(0.14)	462	0
	<i>NoNested</i>	3	171	(0.23)	110	(0.15)	467	0
	<i>ProperShare</i>	9	86	(0.11)	29	(0.04)	633	0
Python	<i>Patience</i>	1	907	(0.21)	561	(0.13)	2 829	1
	<i>NoNested</i>	3	830	(0.19)	602	(0.14)	2 865	1
	<i>ProperShare</i>	9	446	(0.10)	223	(0.05)	3 627	2

TABLE 6.3: Best results for each extraction mode. The height column indicates the minimum height a subtree must have to qualify for sharing, configured with the `--min-height` option. All of the above results were obtained with locally-scoped patches. Globally-scoped success rates were consistently lower than their locally-scoped counterpart.

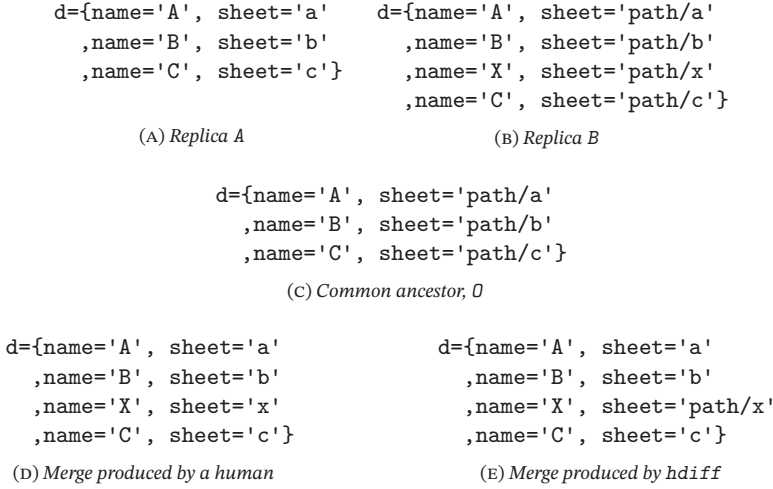


FIGURE 6.2: Example distilled from `hawkthorne-server-lua`, commit 60eba8. One replica introduced entries in a dictionary where another transformed a system path. The `hdiff` tool did produce a correct merge given, but this got classified as *mdif*.

Given the multitude of dials we can adjust in `hdiff`, we have run the experiment with each combination of extraction method (*Patience*, *NoNested*, *ProperShare*), local or global metavariable scoping and minimum sharing height of 1, 3 and 9. Table 6.3 shows the combination of parameters that yielded more successes per extraction method. The column for scoping is omitted because local scope outperformed global scoping in all instances. Table 6.2 shows only the highest success rate per language.

The varying true success rates seen in Table 6.3 are to be expected. Different parameters used with `hdiff` yield different patches, which might be easier or harder to merge. Out of the datapoints that resulted in *mdif* we have manually analyzed 16 randomly selected cases. We witnessed that in 13 of those `hdiff` behaved as we expect, and the *mdif* result was attributed to the human performing more operations than a structural merge would have performed, as exemplified in Figure 6.2, which was distilled from the manually analyzed cases. We will shortly discuss two cases, illustrated in Figures 6.3 and 6.4, where `hdiff` behaved unexpectedly.

It is worth noting that even though 100% success rate is unachievable – some conflicts really come from a subtree being modified in two distinct ways and inevitably require human intervention – the results we have seen are very encouraging. In Table 6.2 we see that `hdiff` produces a merge in at least 39% of datapoints and most often matches the handmade merge.

The cases where *the same* datapoint yields a true success and a *mdiff*, depending on which extraction method was used, are interesting. Let us look at two complementary examples (Figures 6.3 and 6.4) that were distilled from these contradicting cases.

Figure 6.3 shows an example where merging patches extracted with *Patience* returns the correct result, but merging patches extracted with *NoNest* does not. Because replica A modified the definition of `f`, the entire declaration of `f` cannot be copied, and it is placed inside the same scope (alignment) as the definition of `g` since they share a name, `x`. They also share, however, the list of method modifiers, which in this case is `public`. When B modifies the list of modifiers of method `g` by appending `static`, the merge algorithm replicates this change to the list of modifiers of `f`, since the patch wrongly believes both lists represent *the same list*. Merging with *Patience* does not witness the problem since it will not share `x` nor the modifier list, as these occur more than once in the deletion and insertion context of both `hdiff 0 A` and `hdiff 0 B`.

Figure 6.4, on the other hand, shows an example where merging patches extracted with *NoNested* succeeds, but *Patience* inserts a declaration in an unexpected location. Upon further inspection, however, the reason for the diverging behavior becomes clear. When differencing A and 0 under *Patience* context extraction, the empty bodies (which are represented in the Java AST by *MethodBody Nothing*) of the declarations of `n` and `o` are not shared. Hence, the alignment mechanism wrongly identifies that *both* `n` and `o` were deleted. Moreover, because `C.g()` is uniquely shared between the definition of `m` and `S`, the patch identifies that `void m...` became `String S...`. Finally, the merge algorithm then transforms `void m` into `String S`, but then sees two deletions, which trigger the deletion of `n` and `o` from the spine. The next instruction is the insertion of `X`, resulting in the non-intuitive placement of `X` in the merge produced with *Patience*. When using *NoNested*, however, the empty bodies get all shared through the code and prevent the detection of a deletion by the alignment algorithm. It is worth noting that just because Java does not care about the order of declarations, this is not acceptable behavior since it could produce invalid source files in a language like Agda, where the order of declarations matters, for example.

The examples in Figures 6.3 and 6.4 illustrate an inherent difficulty of using naive structured differencing over structures with complex semantics, such as source-code. On the one hand sharing method modifiers triggers undesired replication of a change. On the other, the lack of sharing of empty method bodies makes it difficult to place an insertion in its correct position.

When `hdiff` returned a patch with conflicts, that is, we could *not* successfully solve the merge, we recorded the class of conflicts we observed. Table 6.4 shows the distribution of each conflict type throughout the dataset. Note that a patch resulting from a merge can have multiple conflicts. This information is useful for deciding which aspects of the merge algorithm can yield better results.

```

class Class {
  public int f(int x)
  { F(x*y); }
  public int g(int x)
  { G(x+2); }}
  (A) A.java

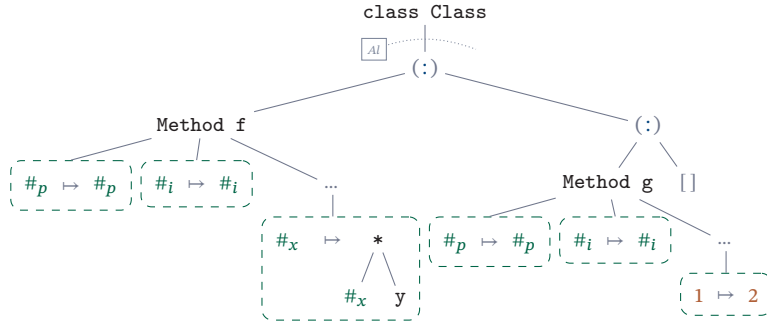
class Class {
  public int f(int x)
  { F(x); }
  public int g(int x)
  { G(x+1); }}
  (B) D.java

class Class {
  public int f(int x)
  { F(x); }
  public static int g(int x)
  { G(x+1); }}
  (C) B.java

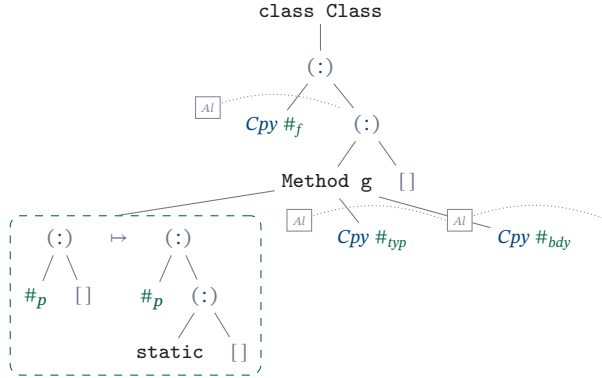
class Class {
  public int f(int x)
  { F(x*y); }
  public static int g(int x)
  { G(x+2); }}
  (D) Expected merge, computed with Patience

class Class {
  public static int f(int x)
  { F(x*y); }
  public static int g(int x)
  { G(x+2); }}
  (E) Incorrect merge, computed with NoNest

```



(F) Simplified illustration of patch computed with `hdiff -d nonest {D,A}.java`. The sharing of $\#_p$ reflects the sharing of the list of method modifiers.



(G) Simplified illustration of patch computed with `hdiff -d nonest {D,B}.java`, note how each copy happens inside its own scope

FIGURE 6.3: Example distilled from `cas`, commit 035eae3, where *Patience* merges with a true success but *NoNest* merges with `mdif`, and, in fact, replicates the `static` modifier incorrectly.

```

class Class {
  String S = C.g();
  void m ()
  { return; }
  void o (int l);
  void p ();
}
(A) A.java

class Class {
  void m ()
  { C.q.g(); return; }
  void n ();
  void o ();
  void p ();
}
(B) D.java

class Class {
  void m ()
  { C.q.g(); return; }
  void n ();
  void o ();
  void X ();
  void p ();
}
(C) B.java

class Class {
  String S = C.g();
  void m ()
  { return; }
  void o (int l);
  void X ();
  void p ();
}
(D) Expected merge, computed with
    NoNested

class Class {
  String S = C.g();
  void X ();
  void m ()
  { return; }
  void o (int l);
  void p ();
}
(E) Incorrect merge, computed
    with Patience

```

FIGURE 6.4: Example distilled from *spring-boot*, commit 0074e9, where *NoNested* merges with a true success but *Patience* merges with *mdiff* since it inserts the declaration of *X* in the wrong place.

	not-eq	inst-mod	del-spn	ins-ins	inst-ins	inst-del	Others
Amount	7904	5052	2144	1892	868	357	506
Ratio	0.42	0.27	0.11	0.1	0.05	0.02	0.03

TABLE 6.4: Distribution of conflicts observed by running *hdiff* over our dataset [72]. The first row displays the number of times that *throwConf* was called with which label.

6.3.1 THREATS TO VALIDITY

The synchronization experiment is encouraging, but before drawing conclusions however, we must analyze our assumptions and setting and preemptively understand which factors could also be influencing the numbers.

We are differencing and comparing objects *after* parsing. This means that comments and formatting data were completely ignored. In fact, preliminary evaluations showed that a vastly inferior success rate results from incorporating and considering source-location tokens in the abstract syntax tree. This is expected since the insertion of a single empty line, for example, will change the hashes that identify all subsequent elements of the abstract syntax and stop them from being shared. The source-location tokens essentially prevent the transformations that happen further down the file to be detected by `hdiff`. Although `stdiff` would not suffer from this problem, it is already impractical by itself.

Our decision to ignore formatting, comments and source-location tokens is twofold. First, the majority of the available parsers does not include said information. Secondly, if we had considered all that information in our merging process, the final numbers would not inform us about how many code transformations are *disjoint* and could be automatically merged.

Another case worth noting is that although we have not found many cases where `hdiff` performed a wrong merge, Figures 6.3 and 6.4 shows two such cases, hence, it is important to take the aggregate success rate with a grain of salt. There exists a probability that some of the *mdif* cases are false positives, that is, `hdiff` produced a merge but it performed the wrong operation.

Finally, one can also argue we have not considered conflicts that arise from rebasing, as these are not observed in the git history. This does not necessarily make a threat to validity, but indeed would have given us more data. That being said, we would only be able to recreate rebases done through the GitHub web interface. The rebases done on the command line are impossible to recreate.

6.4 DISCUSSION

This chapter provided an empirical evaluation of our methods and techniques. We observed how `stdiff` is at least one order of magnitude slower than `hdiff`, confirming our suspicion of it being unusable in practice. Preliminary synchronization experiments done with `stdiff` over the same data revealed a comparatively small success rate. Around 15% of the conflicts could be solved, out of which 60% did match what a human did.

The measurements for `hdiff`, on the other hand, gave impressive results. Even with all the overhead introduced by generic programming and an unoptimized algorithm, we can still compute patches almost instantaneously. Moreover, it confirms our intuition that the differencing algorithm underlying `hdiff` is in fact linear. Moreover, the synchronization results for `hdiff` are encouraging. We have observed that 39% of the conflicts in our dataset could be solved by `hdiff` and 66% of these solutions did match what a human performed.

An interesting observation that comes from the synchronization experiment, Table 6.3, is that the best merging success rate for all languages used the *Patience* context extraction – only copying subtrees that occur uniquely. This suggests that it might be worthwhile to forbid duplication and contractions on the representation level and work on a merging algorithm that enjoys the precondition that each metavariable occurs only twice. This simplification could enable us to write a simpler merging algorithm and an Agda model, which can then be used to prove important properties about our algorithms.



DISCUSSION

Even though the main topic of this thesis is *structural differencing*, a significant part of the contribution lies in the field of generic programming. The two libraries we wrote make it possible to use powerful generic programming techniques over larger classes of datatypes than what was previously available. In particular, defining the generic interpretation as a cofree comonad and a free monad combined in a single datatype is very powerful. Being able to annotate and augment datatypes, for example, was paramount for scaling our algorithms.

On *structural differencing*, we have explored two preliminary approaches. A first method, `stdiff`, was presented in Chapter 4 and revealed itself to be impractical due to poor performance. The second method, `hdiff`, introduced in Chapter 5, has shown much greater potential. Empirical results were discussed in Chapter 6.

7.1 THE FUTURE OF STRUCTURAL DIFFERENCING

The larger picture of structural differencing is more subtle, though. It is not because our preliminary prototype has shown good results that we are ready to scale it to be the next `git merge`. There are three main difficulties in applying structural differencing to source-code with the objective of writing better merge algorithms:

- a) How to properly handle formatting and comments of source code: should the AST keep this information? If so, the tree matching must be adapted to cope with this. Two equal trees must be matched regardless of whether or not they appeared with a different formatting in their respective source files.

- b) How to ensure that subtrees are only being shared within their respective scope and, equally importantly, how to specify which datatypes of the AST are affected by scopes.
- c) When merging fails, returning a patch with conflicts, a human must interact with the tool and solve the conflicts. What kind of interface would be suitable for that? Further ahead, comes the question of automatic conflict solving domain-specific languages. Could we configure the merge algorithm to always chose higher version numbers, for example, whenever it finds a conflict in, say, a config file?

Fixing the obstacles above in a generic way would require a significant effort. So much so that it makes me question the applicability of structural differencing for the exclusive purpose of merging source-code. From a broader perspective, however, there are many other interesting applications that could benefit from structural differencing techniques. In particular, we can probably use structural differencing to aid any task where a human does not directly edit the files being analyzed or when the result of the analysis require no further interaction. For example, it should be possible to deploy `hdiff` to provide a human readable summary of a patch, something that looks at the working directory, computes the structural diffs between the various files, just like `git diff`, but displays information in the lines of:

```
some/project/dir $ hsummary
function fact refactored;
definition of fact changed;
import statements added;
```

In combination with the powerful web interfaces of services like GitHub or GitLab, we could also use tools like `hdiff` to study the evolution of code or to inform the assignee of a pull request whether or not it detected the changes to be *structurally disjoint*. If nothing else, we could at least direct the attention of the developers to the locations in the source-code where there are actual conflicts and the developer has to make a choice. That is where mistakes are more likely to be made. One way of circumventing the formatting and comment issues above is to write a tool that checks whether the developer included all changes in a sensible way and warns them otherwise, but it is always a human performing the actual merge.

Finally, differencing file formats that are based on JSON or XML, such as document processors and spreadsheet processors, might be much easier than source code. Take the formatting of a `.odf` file for example. It is automatically generated and independent of the formatting of document inside the file and it has no scoping or sharing inside, hence, it would be simpler to deploy a structural merging tool over `.odf` files. Some care must be taken with the unordered trees, even though our conjecture is that `hdiff` would behave mostly alright.

7.2 FUTURE WORK

Although the long term future of structural differencing specifically for source-code versioning is uncertain, there are numerous fronts to continue working on many of the aspects developed and discussed in this dissertation, in particular, `hdiff` (Chapter 5). We refer the interested reader to Section 5.4 for a more detailed discussion on these topics, but proceed with a summary of interesting directions for future work.

One clear option for further work is the improvement of the merge algorithm, presented in Section 5.3. A good way to start could be restricting `hdiff` to produce only linear patches (context extraction with *Patience* option) and use these guarantees to study and develop a merge algorithm in a more disciplined fashion. It is possible that the extra guarantees that are provided by linear patches (metavariables are used only once) would simplify the algorithm to the point where we can start thinking about proving properties about it. We would hope that some simplifications would remove the need for some of the more ad-hoc checks that are currently present in the merge algorithm – take the example from Figure 5.32, which feels overly complicated and with no real good justification besides having found these situation in practice. Finally, our experiments have shown us that the *Patience* extraction method gives superior success rates anyway.

Another interesting front would be to define the type of *Chg* in a well-scoped manner, essentially using De Bruijn indicies. This would potentially complicate some of the simpler parts of `hdiff` but could provide important insight into how to handle variables when merging in a very disciplined way. Different representations for *Chg* could also shed some light on how to better control which subtrees can be shared or not.

The actual implementation of `hdiff` could also benefit from further work. We could work on optimizing the generic programming libraries for performance, rewriting parts of the code to use standard implementations well-known data structures instead, or even better visualization of patches using pretty printers.

Finally, the metatheory surrounding `hdiff`’s *Chg* and *Patch* should be worked on. In Section 5.1.3 we have seen how *Chg* forms a partial monoid with a simple composition operation, but we also seen how the trivial inverse operation does not give us a partial group. It could give us an inverse semigroup, for it has a weaker notion of *inverse*. In fact, Darcs patch theory have been formalised with inverse semigroups [47]. Additionally, using the canonical extension order (i.e., comparing domains of application functions) is not a great option for defining *the best* patch. It would be interesting to see whether a categorical approach, similar to Mimram’s work [71], could provide more educated insights in that direction.

7.3 CONCLUDING REMARKS

This dissertation explored a novel approach to structural differencing and a successful prototype for computing and merging patches for said approach. The main novelty comes from relying on unrestricted tree-matchings, which are possible because we never translate to an edit-script-like structure. We have identified the challenges of employing such techniques to merging of source-code but still achieved encouraging empirical results. In the process of developing our prototypes we have also improved the Haskell ecosystem for generic programming.



SOURCE-CODE AND DATASET

A.1 SOURCE-CODE

The easiest way to obtain the source is through either GitHub or Hackage. The source code for the different projects discussed throughout this dissertation is publicly available as Haskell packages, on Hackage:

- `hackage.haskell.org/package/generics-mrsop`
- `hackage.haskell.org/package/simplistic-generics`
- `hackage.haskell.org/package/generics-mrsop-gdiff`
- `hackage.haskell.org/package/hdiff`

The actual version of `hdiff` that we have documented and used to obtain the results presented in this dissertation has been archived on Zenodo [73].

A.2 DATASET

The dataset [72] was obtained by running the data collection script (Section 6.1) over the repositories listed in Table A.1, on the 16th of January of 2020. It is also available in Zenodo for download.

TABLE A.1: *Repositories used for data collection.*

Language	Repository	Conflicts	Commits	Forks
Clojure	metabase/metabase	411	18697	25
Clojure	onyx-platform/onyx	189	6828	209
Clojure	incanter/incanter	96	1593	286
Clojure	nathanmarz/cascalog	68	1366	181
Clojure	overtone/overtone	65	3070	413
Clojure	technomancy/leiningen	46	4736	15
Clojure	ring-clojure/ring	44	1027	441
Clojure	ztellman/aleph	43	1398	213
Clojure	pedestal/pedestal	35	1581	248
Clojure	circleci/frontend	33	18857	170
Clojure	arcadia-unity/Arcadia	25	1716	95
Clojure	walmartlabs/lacinia	19	991	105
Clojure	clojure/clojurescript	18	5706	730
Clojure	oakes/Nightcode	17	1914	119
Clojure	weavejester/compojure	16	943	245
Clojure	boot-clj/boot	12	1331	169
Clojure	clojure-liberator/liberator	12	406	144
Clojure	originrose/cortex	11	1045	103
Clojure	dakrone/clj-http	9	1198	368
Clojure	bhauman/lein-figwheel	9	1833	221
Clojure	jonase/kibit	9	436	124
Clojure	riemann/riemann	7	1717	512
Clojure	korma/Korma	7	491	232
Clojure	clojure/core.async	4	564	181
Clojure	status-im/status-react	3	5224	723
Clojure	cemerick/friend	2	227	122
Clojure	LightTable/LightTable	1	1265	927
Clojure	krisajenkins/yesql	1	285	112
Clojure	cgrand/enlive	1	321	144
Clojure	plumatic/schema	1	825	244
Java	spring-projects/spring-boot	760	24545	284
Java	elastic/elasticsearch	746	49920	158
Java	apereo/cas	363	15834	31
Java	jenkinsci/jenkins	296	29141	6
Java	xetorthio/jedis	147	1610	32
Java	google/ExoPlayer	133	7694	44
Java	apache/storm	117	10204	4
Java	junit-team/junit4	77	2427	29
Java	skylot/jadx	52	1165	24
Java	naver/pinpoint	51	10931	3
Java	apache/beam	34	25062	22
Java	baomidou/mybatis-plus	31	3640	21
Java	mybatis/mybatis-3	21	3164	83
Java	dropwizard/dropwizard	20	5229	31
Java	SeleniumHQ/selenium	18	24627	54

TABLE A.1: *Repositories used for data collection (continued).*

Language	Repository	Conflicts	Commits	Forks
Java	code4craft/webmagic	11	1015	37
Java	aws/aws-sdk-java	7	2340	24
Java	spring-projects/spring-security	7	8339	36
Java	eclipse/deeplearning4j	6	572	48
Java	square/okhttp	5	4407	78
JavaScript	meteor/meteor	1208	22501	51
JavaScript	adobe/brackets	699	17782	66
JavaScript	mrdoob/three.js	403	31473	22
JavaScript	moment/moment	141	3724	65
JavaScript	RocketChat/Rocket.Chat	125	17445	55
JavaScript	serverless/serverless	118	12278	39
JavaScript	nodejs/node	99	29302	159
JavaScript	twbs/bootstrap	86	19261	679
JavaScript	photonstorm/phaser	80	13958	61
JavaScript	emberjs/ember.js	76	19460	42
JavaScript	atom/atom	63	37335	137
JavaScript	TryGhost/Ghost	50	10374	7
JavaScript	jquery/jquery	44	6453	19
JavaScript	mozilla/pdf.js	41	12132	69
JavaScript	Leaflet/Leaflet	37	6810	44
JavaScript	expressjs/express	36	5558	79
JavaScript	hexojs/hexo	27	3146	38
JavaScript	videojs/video.js	17	3509	63
JavaScript	facebook/react	10	12732	273
JavaScript	jashkenas/underscore	8	2447	55
JavaScript	lodash/lodash	8	7992	46
JavaScript	axios/axios	8	900	6
JavaScript	select2/select2	3	2573	58
JavaScript	chartjs/Chart.js	3	2966	101
JavaScript	facebook/jest	2	4595	41
JavaScript	vuejs/vue	1	3076	234
JavaScript	nwjs/nw.js	1	3913	38
Lua	Kong/kong	209	5494	31
Lua	hawkthorne/hawkthorne-journey	155	5538	370
Lua	snabbco/snabb	119	9456	295
Lua	tarantool/tarantool	54	13542	224
Lua	luarocks/luarocks	45	2325	296
Lua	luakit/luakit	28	4186	219
Lua	pkulchenko/ZeroBraneStudio	20	3945	447
Lua	CorsixTH/CorsixTH	16	3355	250
Lua	OpenNMT/OpenNMT	14	1684	455
Lua	koreader/koreader	14	7256	710
Lua	bakpakin/Fennel	12	689	59
Lua	Olivine-Labs/busted	9	950	139

TABLE A.1: *Repositories used for data collection (continued).*

Language	Repository	Conflicts	Commits	Forks
Lua	Element-Research/rnn	8	622	318
Lua	lcpz/awesome-copypcats	8	821	412
Lua	Tieske/Penlight	6	743	190
Lua	yagop/telegram-bot	5	729	519
Lua	awesomeWM/awesome	5	9990	360
Lua	torch/nn	4	1839	967
Lua	luvit/luvit	4	2897	330
Lua	GUI/lua-resty-auto-ssl	3	318	119
Lua	alexazhou/VeryNginx	3	604	810
Lua	sailorproject/sailor	2	640	128
Lua	leafo/moonscript	2	738	162
Lua	nrk/redis-lua	1	327	193
Lua	skywind3000/z.lua	1	367	59
Lua	rxl/json.lua	1	46	144
Lua	luafun/luafun	1	55	88
Python	python/cpython	891	106167	131
Python	sympy/sympy	864	41009	29
Python	matplotlib/matplotlib	515	32949	47
Python	home-assistant/home-assistant	496	23812	91
Python	bokeh/bokeh	326	18196	32
Python	certbot/certbot	272	9524	28
Python	scikit-learn/scikit-learn	192	25044	19
Python	explosion/spaCy	163	11141	27
Python	docker/compose	129	5590	29
Python	scrapy/scrapy	74	7705	83
Python	keras-team/keras	70	5342	176
Python	tornadoweb/tornado	60	4144	51
Python	pallets/flask	56	3799	132
Python	ipython/ipython	51	24203	39
Python	pandas-dev/pandas	48	21596	92
Python	quantopian/zipline	45	6032	31
Python	Theano/Theano	44	28099	25
Python	psf/requests	32	5927	75
Python	ansible/ansible	29	48864	18
Python	nvbn/thefuck	11	1555	26
Python	waditu/tushare	8	407	35
Python	facebook/prophet	4	445	26
Python	jakubroztocil/httpie	3	1145	29
Python	binux/pyspider	1	1174	34
Python	Jack-Cherish/python-spider	1	279	39
Python	zulip/zulip	1	34149	35

BIBLIOGRAPHY

- [1] ADAMS, M. D. Scrap your zippers: A generic zipper for heterogeneous types. In *WGP '10: Proceedings of the 2010 ACM SIGPLAN workshop on Generic programming* (New York, NY, USA, 2010), ACM, pp. 13–24.
- [2] AKUTSU, T. Tree edit distance problems: Algorithms and applications to bioinformatics. *IEICE Transactions on Information and Systems E93.D*, 2 (2010), 208–218.
- [3] AKUTSU, T., FUKAGAWA, D., AND TAKASU, A. Approximating tree edit distance through string edit distance. *Algorithmica* 57, 2 (Jun 2010), 325–348.
- [4] ALIMARINE, A. *Generic Functional Programming – Conceptual Design, Implementation and Applications*. PhD thesis, Radboud University, 2005.
- [5] ALTENKIRCH, T., GHANI, N., HANCOCK, P., MCBRIDE, C., AND MORRIS, P. Indexed containers. *Journal of Functional Programming* 25 (2015).
- [6] ANDERSEN, J., AND LAWALL, J. L. Generic patch inference. In *23rd IEEE/ACM International Conference on Automated Software Engineering* (L'Aquila, Italy, Sept. 2008), pp. 337–346.
- [7] ASENOV, D., GUENAT, B., MÜLLER, P., AND OTTH, M. 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* (New York, NY, USA, 2017), Springer-Verlag New York, Inc., pp. 152–169.
- [8] AUGSTEN, N., BOHLEN, M., DYRESON, C., AND GAMPER, J. Approximate joins for data-centric XML. In *2008 IEEE 24th International Conference on Data Engineering* (2008), IEEE, pp. 814–823.
- [9] AUGSTEN, N., BÖHLEN, M., AND GAMPER, J. The pq-gram distance between ordered labeled trees. *ACM Transactions on Database Systems (TODS)* 35, 1 (2010), 4.
- [10] AUTEXIER, S. Similarity-based diff, three-way diff and merge. *Int. J. Software and Informatics* 9, 2 (2015), 259–277.
- [11] BACKURS, A., AND INDYK, P. Edit distance cannot be computed in strongly sub-quadratic time (unless seth is false). In *Proceedings of the Forty-Seventh Annual ACM Symposium on Theory of Computing* (New York, NY, USA, 2015), STOC '15, Association for Computing Machinery, p. 51–58.

- [12] BALASUBRAMANIAM, S., AND PIERCE, B. C. What is a file synchronizer? In *Proceedings of the 4th Annual ACM/IEEE International Conference on Mobile Computing and Networking* (New York, NY, USA, 1998), MobiCom '98, ACM, pp. 98–108.
- [13] BARENDREGT, H. P. *The Lambda Calculus: Its Syntax and Semantics*. Elsevier, 1984.
- [14] BERGROTH, L., HAKONEN, H., AND RAITA, T. A survey of longest common subsequence algorithms. In *Proceedings of the Seventh International Symposium on String Processing Information Retrieval (SPIRE'00)* (Washington, DC, USA, 2000), SPIRE '00, IEEE Computer Society, pp. 39–.
- [15] BEZEM, M., KLOP, J., BARENDSEN, E., DE VRIJER, R., AND TERESE. *Term Rewriting Systems*. Cambridge Tracts in Theoretical Computer Science. Cambridge University Press, 2003.
- [16] BILLE, P. A survey on tree edit distance and related problems. *Theor. Comput. Sci.* 337, 1-3 (June 2005), 217–239.
- [17] BIRD, R., AND MEERTENS, L. Nested datatypes. In *International Conference on Mathematics of Program Construction* (1998), Springer, pp. 52–67.
- [18] BOTTU, G.-J., KARACHALIAS, G., SCHRIJVERS, T., OLIVEIRA, B. C. D. S., AND WADLER, P. Quantified class constraints. In *Proceedings of the 10th ACM SIGPLAN International Symposium on Haskell* (New York, NY, USA, 2017), Haskell 2017, ACM, pp. 148–161.
- [19] BRASS, P. *Advanced Data Structures*. Cambridge University Press, 2008.
- [20] BRAVENBOER, M., KALLEBERG, K. T., VERMAAS, R., AND VISSER, E. Stratego/xt 0.17. a language and toolset for program transformation. *Science of computer programming* 72, 1-2 (2008), 52–70.
- [21] CHAWATHE, S. S., AND GARCIA-MOLINA, H. Meaningful change detection in structured data. In *Proceedings of the 1997 ACM SIGMOD International Conference on Management of Data* (New York, NY, USA, 1997), SIGMOD '97, ACM, pp. 26–37.
- [22] CHAWATHE, S. S., RAJARAMAN, A., GARCIA-MOLINA, H., AND WIDOM, J. Change detection in hierarchically structured information. In *Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data* (New York, NY, USA, 1996), SIGMOD '96, ACM, pp. 493–504.
- [23] COBÉNA, G., ABITEBOUL, S., AND MARIAN, A. Detecting changes in XML documents. In *Proceedings - International Conference on Data Engineering* (02 2002), pp. 41–52.

- [24] COHEN, B. Patience diff advantages, 2010. <https://bramcohen.livejournal.com/73318.html>.
- [25] DAMERAU, F. J. A technique for computer detection and correction of spelling errors. *Commun. ACM* 7, 3 (Mar. 1964), 171–176.
- [26] DE VRIES, E., AND LÖH, A. True sums of products. In *Proceedings of the 10th ACM SIGPLAN Workshop on Generic Programming* (New York, NY, USA, 2014), WGP '14, ACM, pp. 83–94.
- [27] DEMAINE, E. D., MOZES, S., ROSSMAN, B., AND WEIMANN, O. An optimal decomposition algorithm for tree edit distance. In *Proceedings of the 34th International Colloquium on Automata, Languages and Programming (ICALP 2007)* (Wroclaw, Poland, July 9–13 2007), pp. 146–157.
- [28] DULUCQ, S., AND TOUZET, H. Analysis of tree edit distance algorithms. In *Combinatorial Pattern Matching* (Berlin, Heidelberg, 2003), R. Baeza-Yates, E. Chávez, and M. Crochemore, Eds., Springer Berlin Heidelberg, pp. 83–95.
- [29] EISENBERG, R. A., AND WEIRICH, S. Dependently typed programming with singletons. In *Proceedings of the 2012 Haskell Symposium* (New York, NY, USA, 2012), Haskell '12, ACM, pp. 117–130.
- [30] EISENBERG, R. A., WEIRICH, S., AND AHMED, H. G. Visible type application. In *Programming Languages and Systems - 25th European Symposium on Programming, ESOP 2016, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2016, Eindhoven, The Netherlands, April 2-8, 2016, Proceedings* (2016), pp. 229–254.
- [31] FALLERI, J., MORANDAT, F., BLANC, X., MARTINEZ, M., AND MONPERRUS, M. Fine-grained and accurate source code differencing. In *ACM/IEEE International Conference on Automated Software Engineering, ASE '14, Vasteras, Sweden - September 15 - 19, 2014* (2014), pp. 313–324.
- [32] FARINIER, B., GAZAGNAIRE, T., AND MADHAVAPEDDY, A. Mergeable persistent data structures. In *Vingt-sixièmes Journées Francophones des Langages Applicatifs (JFLA 2015)* (Le Val d'Ajol, France, Jan. 2015), D. Baelde and J. Alglave, Eds.
- [33] FILLIÂTRE, J.-C., AND CONCHON, S. Type-safe modular hash-consing. In *Proceedings of the 2006 Workshop on ML* (New York, NY, USA, 2006), ML '06, ACM, pp. 12–19.
- [34] FINIS, J. P., RAIBER, M., AUGSTEN, N., BRUNEL, R., KEMPER, A., AND FÄRBER, F. Rws-diff: Flexible and efficient change detection in hierarchical data. In *Proceedings of the 22nd ACM International Conference on Information & Knowledge Management* (New York, NY, USA, 2013), CIKM '13, ACM, pp. 339–348.

- [35] FOSTER, J. N., GREENWALD, M. B., KIRKEGAARD, C., PIERCE, B. C., AND SCHMITT, A. Exploiting schemas in data synchronization. *Journal of Computer and System Sciences* 73, 4 (2007), 669–689.
- [36] GARUFFI, G. Version control systems: Diffing with structure. Master's thesis, Utrecht Universiteit, 2018.
- [37] GHANI, N., LÜTH, C., DE MARCHI, F., AND POWER, J. Algebras, coalgebras, monads and comonads. *Electronic Notes in Theoretical Computer Science* 44, 1 (2001), 128–145.
- [38] GIBBONS, J. Design patterns as higher-order datatype-generic programs. In *Proceedings of the 2006 ACM SIGPLAN Workshop on Generic Programming* (New York, NY, USA, 2006), WGP '06, ACM, pp. 1–12.
- [39] GOTO, E. Monocopy and associative algorithms in an extended lisp. Tech. rep., University of Tokyo, 1974.
- [40] GUHA, S., JAGADISH, H. V., KOUDAS, N., SRIVASTAVA, D., AND YU, T. Approximate XML joins. In *Proceedings of the 2002 ACM SIGMOD International Conference on Management of Data* (New York, NY, USA, 2002), SIGMOD '02, Association for Computing Machinery, p. 287–298.
- [41] HASHIMOTO, M., AND MORI, A. Diff/ts: A tool for fine-grained structural change analysis. In *Proceedings of the 2008 15th Working Conference on Reverse Engineering* (Washington, DC, USA, 2008), WCRE '08, IEEE Computer Society, pp. 279–288.
- [42] HENIKOFF, S., AND HENIKOFF, J. G. Amino acid substitution matrices from protein blocks. *Proceedings of the National Academy of Sciences* 89, 22 (1992), 10915–10919.
- [43] HINZE, R., JEURING, J., AND LÖH, A. Type-indexed data types. *Science of Computer Programming* 51, 1 (2004), 117 – 151. Mathematics of Program Construction (MPC 2002).
- [44] HIRSCHBERG, D. S. A linear space algorithm for computing maximal common subsequences. *Commun. ACM* 18, 6 (June 1975), 341–343.
- [45] HUET, G. The zipper. *J. Funct. Program.* 7 (09 1997), 549–554.
- [46] HUNT, J. W., AND MCILROY, M. D. An algorithm for differential file comparison. Tech. Rep. CSTR 41, Bell Laboratories, Murray Hill, NJ, 1976.
- [47] JACOBSON, J. A formalization of darcs patch theory using inverse semigroups. Tech. rep., UCLA, 2009.

- [48] KHANNA, S., KUNAL, K., AND PIERCE, B. C. A formal investigation of diff3. In *FSTTCS 2007: Foundations of Software Technology and Theoretical Computer Science* (Berlin, Heidelberg, 2007), V. Arvind and S. Prasad, Eds., Springer Berlin Heidelberg, pp. 485–496.
- [49] KIM, D., NAM, J., SONG, J., AND KIM, S. Automatic patch generation learned from human-written patches. In *Proceedings of the 2013 International Conference on Software Engineering* (2013), IEEE Press, pp. 802–811.
- [50] KLEIN, P. N. Computing the edit-distance between unrooted ordered trees. In *Proceedings of the 6th Annual European Symposium on Algorithms* (London, UK, UK, 1998), ESA '98, Springer-Verlag, pp. 91–102.
- [51] KLINT, P., VAN DER STORM, T., AND VINJU, J. Rascal: A domain specific language for source code analysis and manipulation. In *2009 Ninth IEEE International Working Conference on Source Code Analysis and Manipulation* (2009), IEEE, pp. 168–177.
- [52] KNUTH, D. E. The genesis of attribute grammars. In *Proceedings of the International Conference on Attribute Grammars and Their Applications* (Berlin, Heidelberg, 1990), WAGA, Springer-Verlag, p. 1–12.
- [53] LÄMMEL, R., AND PEYTON JONES, S. Scrap your boilerplate: A practical design pattern for generic programming. In *Proceedings of the 2003 ACM SIGPLAN International Workshop on Types in Languages Design and Implementation* (New York, NY, USA, 2003), TLDI '03, ACM, pp. 26–37.
- [54] LANHAM, M., KANG, A., HAMMER, J., HELAL, A., AND WILSON, J. Format-independent change detection and propagation in support of mobile computing. *Brazilian Symposium on Databases (SBBD)* (01 2002), 27–41.
- [55] LEMPSINK, E., LEATHER, S., AND LÖH, A. Type-safe diff for families of datatypes. In *Proceedings of the 2009 ACM SIGPLAN Workshop on Generic Programming* (New York, NY, USA, 2009), WGP '09, ACM, pp. 61–72.
- [56] LEVENSHTAIN, V. Binary codes capable of correcting deletions, insertions and reversals. *Soviet Physics Doklady* 10 (1966), 707.
- [57] LINDHOLM, T. A three-way merge for XML documents. In *Proceedings of the 2004 ACM Symposium on Document Engineering* (New York, NY, USA, 2004), DocEng '04, ACM, pp. 1–10.
- [58] LÖH, A., AND MAGALHAES, J. P. Generic programming with indexed functors. In *Proceedings of the seventh ACM SIGPLAN workshop on Generic programming* (2011), ACM, pp. 1–12.

- [59] MAGALHÃES, J. P., DIJKSTRA, A., JEURING, J., AND LÖH, A. A generic deriving mechanism for Haskell. In *Proceedings of the Third ACM Haskell Symposium on Haskell* (New York, NY, USA, 2010), Haskell '10, ACM, pp. 37–48.
- [60] MAGALHÃES, J. P., AND LÖH, A. A formal comparison of approaches to datatype-generic programming. In *Proceedings Fourth Workshop on Mathematically Structured Functional Programming, Tallinn, Estonia, 25 March 2012* (2012), J. Chapman and P. B. Levy, Eds., vol. 76 of *Electronic Proceedings in Theoretical Computer Science*, Open Publishing Association, pp. 50–67.
- [61] MAGALHÃES, J. P., AND LÖH, A. Generic generic programming. In *International Symposium on Practical Aspects of Declarative Languages* (2014), Springer, pp. 216–231.
- [62] MALETIC, J. I., AND COLLARD, M. L. Exploration, analysis, and manipulation of source code using srcml. In *Proceedings of the 37th International Conference on Software Engineering - Volume 2* (2015), ICSE '15, IEEE Press, p. 951–952.
- [63] MARLOW, S., ET AL. Haskell 2010 Language Report, 2010. <https://www.haskell.org/onlinereport/haskell2010/>.
- [64] MCBRIDE, C. The derivative of a regular type is its type of one-hole contexts (extended abstract), 2001.
- [65] MCBRIDE, C. *Epigram: Practical Programming with Dependent Types*. Springer Berlin Heidelberg, Berlin, Heidelberg, 2005, pp. 130–170.
- [66] MCBRIDE, C., AND MCKINNA, J. The view from the left. *J. Funct. Program.* 14, 1 (Jan. 2004), 69–111.
- [67] MCKENNA, A., HANNA, M., BANKS, E., SIVACHENKO, A., CIBULSKIS, K., KERNYTSKY, A., GARIMELLA, K., ALTSHULER, D., GABRIEL, S., DALY, M., AND DEPRISTO, M. The genome analysis toolkit: A mapreduce framework for analyzing next-generation dna sequencing data. *Genome research* 20 (09 2010), 1297–303.
- [68] MEDEIROS, F., RIBEIRO, M., GHEYI, R., APEL, S., KÄSTNER, C., FERREIRA, B., CARVALHO, L., AND FONSECA, B. Discipline matters: Refactoring of preprocessor directives in the# ifdef hell. *IEEE Transactions on Software Engineering* 44, 5 (2017), 453–469.
- [69] MENEZES A. J., P. V. O., AND VANSTONE, S. A. *Handbook of Applied Cryptography*, boca raton, xiii, 780 ed. CRC Press, 1997.
- [70] MERKLE, R. C. A digital signature based on a conventional encryption function. In *Advances in Cryptology — CRYPTO '87* (Berlin, Heidelberg, 1988), C. Pomerance, Ed., Springer Berlin Heidelberg, pp. 369–378.

- [71] MIMRAM, S., AND GIUSTO, C. D. A categorical theory of patches. *CoRR abs/1311.3903* (2013).
- [72] MIRALDO, V. C. Dataset of merge conflicts collected from github repositories. <https://doi.org/10.5281/zenodo.3751038>, Apr. 2020.
- [73] MIRALDO, V. C. hdiff: hash-based differencing of structured data. <https://doi.org/10.5281/zenodo.3754632>, Apr. 2020.
- [74] MIRALDO, V. C., DAGAND, P.-E., AND SWIERSTRA, W. Type-directed diffing of structured data. In *Proceedings of the 2Nd ACM SIGPLAN International Workshop on Type-Driven Development* (New York, NY, USA, 2017), TyDe 2017, ACM, pp. 2–15.
- [75] MIRALDO, V. C., AND SERRANO, A. Sums of products for mutually recursive datatypes: the appropriationist’s view on generic programming. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Type-Driven Development* (2018), ACM, pp. 65–77.
- [76] MIRALDO, V. C., AND SWIERSTRA, W. An efficient algorithm for type-safe structural diffing. *PACMPL* 3, ICFP (2019), 113:1–113:29.
- [77] MITCHELL, N., AND RUNCIMAN, C. Uniform boilerplate and list processing. In *Proceedings of the ACM SIGPLAN Workshop on Haskell Workshop* (New York, NY, USA, 2007), Haskell ’07, ACM, pp. 49–60.
- [78] MOUAT, A. XML diff and patch utilities. Master’s thesis, Heriot-Watt University, Edinburgh, 2002.
- [79] MYERS, E. W. An $O(nd)$ difference algorithm and its variations. *Algorithmica* 1 (1986), 251–266.
- [80] NAVARRO, G. A guided tour to approximate string matching. *ACM Comput. Surv.* 33, 1 (Mar. 2001), 31–88.
- [81] NERON, P., TOLMACH, A., VISSER, E., AND WACHSMUTH, G. A theory of name resolution. In *Programming Languages and Systems* (Berlin, Heidelberg, 2015), J. Vitek, Ed., Springer Berlin Heidelberg, pp. 205–231.
- [82] NOORT, T. V., RODRIGUEZ, A., HOLDERMANS, S., JEURING, J., AND HEEREN, B. A lightweight approach to datatype-generic rewriting. In *Proceedings of the ACM SIGPLAN Workshop on Generic Programming* (New York, NY, USA, 2008), WGP ’08, ACM, pp. 13–24.
- [83] NORELL, U. Dependently typed programming in agda. In *International school on advanced functional programming* (2008), Springer, pp. 230–266.

- [84] PAAßEN, B. Revisiting the tree edit distance and its backtracing: A tutorial. *CoRR abs/1805.06869* (2018).
- [85] PAAßEN, B., HAMMER, B., PRICE, T. W., BARNES, T., GROSS, S., AND PINKWART, N. The continuous hint factory - providing hints in vast and sparsely populated edit distance spaces. *CoRR abs/1708.06564* (2017).
- [86] PALIX, N., THOMAS, G., SAHA, S., CALVÈS, C., LAWALL, J., AND MULLER, G. Faults in Linux: Ten years later. In *Sixteenth International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS 2011)* (Newport Beach, CA, USA, Mar. 2011).
- [87] PAWLIK, M., AND AUGSTEN, N. RTED: A robust algorithm for the tree edit distance. *CoRR abs/1201.0230* (2012).
- [88] PETERS, L. Change detection in XML trees: a survey. In *3rd Twente Student Conference on IT* (2005).
- [89] PEYTON JONES, S., WEIRICH, S., EISENBERG, R. A., AND VYTINIOTIS, D. A reflection on types. In *A List of Successes That Can Change the World*. Springer, 2016, pp. 292–317.
- [90] PICKERING, M., ÉRDI, G., PEYTON JONES, S., AND EISENBERG, R. A. Pattern synonyms. In *Proceedings of the 9th International Symposium on Haskell* (New York, NY, USA, 2016), Haskell 2016, ACM, pp. 80–91.
- [91] PLOTKIN, G. Further note on inductive generalization. *Machine Intelligence*. 6 (01 1971).
- [92] PUTTEN, A. v. Generic diffing and merging of mutually recursive datatypes in Haskell. Master's thesis, Utrecht Universiteit, 2019.
- [93] ROBINSON, E., AND ROSOLINI, G. Categories of partial maps. *Information and computation* 79, 2 (1988), 95–130.
- [94] ROBINSON, J. A. A machine-oriented logic based on the resolution principle. *J. ACM* 12, 1 (Jan. 1965), 23–41.
- [95] RODRIGUEZ, A., JEURING, J., JANSSON, P., GERDES, A., KISELYOV, O., AND OLIVEIRA, B. C. D. S. Comparing libraries for generic programming in Haskell. In *Proceedings of the First ACM SIGPLAN Symposium on Haskell* (New York, NY, USA, 2008), Haskell '08, ACM, pp. 111–122.
- [96] ROY CHOUDHURY, R., YIN, H., AND FOX, A. Scale-driven automatic hint generation for coding style. In *Intelligent Tutoring Systems* (Cham, 2016), A. Micarelli, J. Stamper, and K. Panourgia, Eds., Springer International Publishing, pp. 122–132.

- [97] SERRANO, A., AND HAGE, J. Generic matching of tree regular expressions over Haskell data types. In *Practical Aspects of Declarative Languages - 18th International Symposium, PADL 2016, St. Petersburg, FL, USA, January 18-19, 2016. Proceedings* (2016), pp. 83–98.
- [98] SERRANO, A., AND MIRALDO, V. C. Generic programming of all kinds. In *Proceedings of the 11th ACM SIGPLAN International Symposium on Haskell* (2018), ACM, pp. 41–54.
- [99] SHAPIRA, D., AND STORER, J. A. Edit distance with move operations. In *Combinatorial Pattern Matching* (Berlin, Heidelberg, 2002), A. Apostolico and M. Takeda, Eds., Springer Berlin Heidelberg, pp. 85–98.
- [100] SHEARD, T., AND PEYTON JONES, S. Template meta-programming for Haskell. In *Proceedings of the 2002 ACM SIGPLAN Workshop on Haskell* (New York, NY, USA, 2002), Haskell '02, ACM, pp. 1–16.
- [101] SMITH, R. Version 3.7, December 2018; distributed with GNU diffutils package, 1988. <http://ftp.gnu.org/gnu/diffutils/>.
- [102] TAI, K.-C. The tree-to-tree correction problem. *J. ACM* 26, 3 (July 1979), 422–433.
- [103] THIJE, J., AND ZEEVAERT, L. *Receptive Multilingualism: Linguistic Analyses, Language Policies, and Didactic Concepts*. Hamburg studies on multilingualism. J. Benjamins Publishing Company, 2007.
- [104] TILMANN, F. latexdiff, 2004. mirrors.ctan.org/support/latexdiff/doc/latexdiff-man.pdf.
- [105] VAN HOREBEEK, I., AND LEWI, J. *Abstract Data Types as Initial Algebras*. Springer Berlin Heidelberg, Berlin, Heidelberg, 1989, pp. 14–65.
- [106] VAN TONDER, R., AND LE GOUES, C. Lightweight multi-language syntax transformation with parser parser combinators. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation* (2019), ACM, pp. 363–378.
- [107] VASSENA, M. Generic diff3 for algebraic datatypes. In *Proceedings of the 1st International Workshop on Type-Driven Development, TyDe@ICFP 2016, Nara, Japan, September 18, 2016* (2016), pp. 62–71.
- [108] WADLER, P. Views: A way for pattern matching to cohabit with data abstraction. In *Proceedings of the 14th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages* (1987), POPL '87, pp. 307–313.

- [109] WEIRICH, S., HSU, J., AND EISENBERG, R. A. System FC with explicit kind equality. In *Proceedings of the 18th ACM SIGPLAN International Conference on Functional Programming* (New York, NY, USA, 2013), ICFP '13, Association for Computing Machinery, p. 275–286.
- [110] WEIRICH, S., VOIZARD, A., DE AMORIM, P. H. A., AND EISENBERG, R. A. A specification for dependent types in Haskell. *Proc. ACM Program. Lang.* 1, ICFP (Aug. 2017), 31:1–31:29.
- [111] XI, H., CHEN, C., AND CHEN, G. Guarded recursive datatype constructors. In *Proceedings of the 30th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (New York, NY, USA, 2003), POPL '03, ACM, pp. 224–235.
- [112] YAKUSHEV, A. R., HOLDERMANS, S., LÖH, A., AND JEURING, J. Generic programming with fixed points for mutually recursive datatypes. In *Proceedings of the 14th ACM SIGPLAN International Conference on Functional Programming* (New York, NY, USA, 2009), ICFP '09, ACM, pp. 233–244.
- [113] YORGEY, B. A., WEIRICH, S., CRETIN, J., PEYTON JONES, S., VYTINIOTIS, D., AND MAGALHÃES, J. P. Giving Haskell a promotion. In *Proceedings of the 8th ACM SIGPLAN Workshop on Types in Language Design and Implementation* (New York, NY, USA, 2012), TLDI '12, ACM, pp. 53–66.
- [114] ZHANG, K., AND SHASHA, D. Simple fast algorithms for the editing distance between trees and related problems. *SIAM journal on computing* 18, 6 (1989), 1245–1262.
- [115] ZHANG, K., STATMAN, R., AND SHASHA, D. On the editing distance between unordered labeled trees. *Information processing letters* 42, 3 (1992), 133–139.

SAMENVATTING

De UNIX `diff` tool – die het verschil tussen twee bestanden berekent in termen van de verzameling regels die gekopieerd worden – wordt veel gebruikt bij het versiebeheer van software. De vaste granulariteit, *regels code*, is helaas soms te grof en verhult eenvoudige wijzigingen. Bijvoorbeeld, door het hernoemen van een parameter van een functie kunnen vele verschillende regels veranderen. Dit kan leiden tot onnodige conflicten wanneer ongerelateerde wijzigingen toevallig op dezelfde regel code plaatsvinden. Het is daarom lastig om zulke wijzigingen automatisch te combineren.

Traditionele methodes om verschillen te bepalen tussen een bronbestand en doelbestand [16, 14, 84, 46, 79, 44] maken gebruik van een *edit-script*, die de veranderingen documenteren om het bronbestand in het doelbestand te transformeren. Zulke edit-scripts bestaan uit een reeks edit-operaties, zoals het invoegen, verwijderen of kopiëren van een regel. Beschouw bijvoorbeeld de volgende twee bestanden:

<pre>1 res := 0; 2 for (i in is) { 3 res += i; 4 }</pre>	<pre>1 print("summing up"); 2 sum := 0; 3 for (i in is) { 4 sum += i; 5 }</pre>
--	--

Regels 2 en 4 in het bronbestand links komen overeen met regels 3 en 5 in het doelbestand rechts. Deze worden dan ook geïdentificeerd als kopieën. De overige regels worden verwijderd of ingevoegd. In dit voorbeeld worden regels 1 en 3 uit het bronbestand verwijderd; regels 1,2 en 4 worden in het doelbestand ingevoegd.

Deze informatie over welke afzonderlijke regels zijn gekopieerd, verwijderd of ingevoegd wordt dan samengebracht in een edit script: een lijst operaties die een bronbestand transformeert in een doelbestand. In het voorbeeld hierboven, zou het edit-script bestaan uit een serie edit-operaties als: verwijder een regel; voeg twee nieuwe regels in; kopieer een regel; verwijder een regel; enz. De uitvoer van UNIX `diff` bestaat uit zo'n lijst operaties. Verwijderingen worden aangeduid door een regel te beginnen met een minteken; invoegingen worden aangeduid met een plusteken. In ons voorbeeld zou het resultaat van UNIX `diff` bestaan uit de volgende regels:

```
-  res := 0;
+  print("summing up");
+  sum := 0;
+  for (i in is) {
-    res += i;
+    sum += i;
+  }
```

Er bestaan veel generalisaties van edit-scripts die niet werken met regels code, maar bomen [114, 27, 28, 87, 8, 9], maar veel van dit werk heeft significante nadelen. Om te

beginnen, edit-scripts zijn niet in staat om willekeurige permutaties, duplicaties, of contracties (de inverse van duplicaties) uit te drukken. Ten tweede, hebben de meeste van deze algoritmen een significant slechtere *performance* dan `UNIX diff`. Tot slot, houdt het meeste van dit werk zich bezig met *ongetypeerde* bomen, dat wil zeggen, dat ze geen garanties geven over of de resulterende edit-scripts een zekere structuur, ook wel bekend als een *schema*, behoudt. Het is mogelijk om getypeerde edit-scripts te ontwerpen [55], maar dit pakt nog niet alle bovengenoemde nadelen aan.

In dit proefschrift bespreken we twee nieuwe manieren om het verschil te bepalen tussen ‘gestructureerde data’, zoals bomen, die zich niet langer beperken tot alleen edit-scripts. De eerste aanpak definieert een type-geïndexeerde representatie van wijzigingen en geeft een helder algoritme om twee verschillende wijzigingen samen te voegen, maar is helaas computationeel te duur om bruikbaar te zijn. De tweede aanpak is een stuk efficiënter door het kiezen van een meer extensionele representatie van wijzigingen. Hierdoor kunnen we allerlei transformaties uitdrukken die gebruik maken van invoegingen, verwijderingen, duplicaties, contracties en permutaties, en deze in lineaire tijd berekenen.

Stel, we moeten een wijziging beschrijven in de linkerdeelboom van een binaire boom. Als we een hele programmeertaal zoals Haskell tot onze beschikking zouden hebben, zouden we iets kunnen schrijven als de functie *c* in Figure A.1(A). Merk hierbij op dat deze functie een duidelijk domein heeft – de verzameling bomen die, wanneer *c* erop toegepast wordt een *Just* constructor opleveren. Dit domein wordt bepaald door de patronen en “guards” die de functie *c* gebruikt. Zo bepalen we voor elke boom in dit domein, een bijbehorend resultaat in het codomein. Deze nieuwe boom kunnen we bepalen door in de oude boom de waarde 10 te vervangen door 42. Bij nadere inspectie, zien we dat we het pattern-matchen op de invoerboom kunnen opvatten als het (mogelijk) *verwijderen* van bepaalde deelbomen; de constructie van de resultaatboom kunnen we beschouwen als het *invoegen* van nieuwe deelbomen. Het `hdiff` algoritme dat wij in dit proefschrift ontwikkelen representeert een wijziging *c* dan ook als een patroon en een expressie. Zo kunnen we de wijziging van *c* beschrijven als *Chg (Bin (Leaf 10) y) (Bin (Leaf 42) y)* – zoals we illustreren in Figure A.1(B).

Doordat onze wijzigingen een rijkere expressiviteit genieten, wordt het samenvoegen van zulke wijzigingen complexer. Als gevolg hiervan wordt het algoritme om twee wijzigingen te verenigen ingewikkelder en kan het soms lastiger worden om over de wijzigingen te redeneren.

Deze twee verschillende algoritmes om het verschil tussen gestructureerde data te berekenen kunnen worden toegepast op wederzijds recursieve datatypes, met als gevolg dat ze gebruikt kunnen worden om computerprogramma’s te vergelijken. Om dit te implementeren was niet eenvoudig, en we hebben in de context van dit proefschrift dan ook verschillende nieuwe bibliotheken voor generiek programmeren in Haskell ontwikkeld.

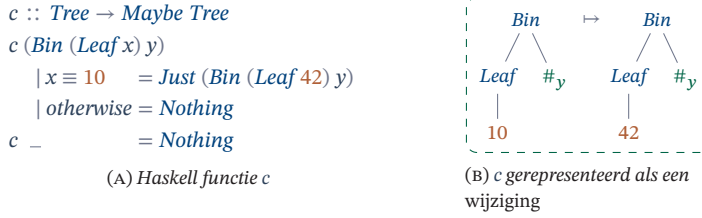


FIGURE A.1: Een Haskell functie en de bijbehorende grafische representatie als wijziging. Deze wijziging past de linkerdeelboom van een binaire knoop aan. De notatie $\#_y$ wordt gebruikt om aan te geven dat y een metavariable is.

Tot slot, hebben we een empirische evaluatie van onze algoritmes uitgevoerd aan de hand van conflicten die we hebben verzameld van GitHub. Uit deze evaluatie blijkt dat ten minste 26% van de conflicten die softwareontwikkelaars dagelijks tegenkomen, voorkomen hadden kunnen worden met de technologie die in dit proefschrift ontwikkeld wordt. Dit doet vermoeden dat er nog veel winst te behalen is in het gebruik van betere algoritmes als basis voor het versiebeheer van software.

CURRICULUM VITAE

Victor Cacciari Miraldo

Born 16 October 1991 in São Paulo, Brazil.

EDUCATION *2009 - 2012, Universidade do Minho, Portugal*
Bachelor in Computer Science
2012 - 2015, Universidade do Minho, Portugal
Master in Computer Science
2015 - 2020, Utrecht Universiteit, the Netherlands
PhD in Computer Science

RESEARCH *2011 - 2012, Universidade do Minho, Portugal*
Introduction to Research Scholarship
2013 - 2014, Universidade do Minho, Portugal
Research Assistant (AVIACC project)
2014 - 2015, Universidade do Minho, Portugal
Research Assistant (QAIS project)
2017 - 2017, Oracle Labs, USA
Internship
2017 - 2020, Oracle Labs, USA (remote)
Research Assistant (Distributed Systems Group)

ACKNOWLEDGEMENTS

This thesis have been reviewed by Jeremy Gibbons, José Nuno Oliveira, Rinus Plasmeijer, Sven-Bodo Scholz, Johan Jeuring; I would like to thank them for their time and for providing careful feedback. I will be defending this thesis against the aforementioned reading committee, Arno Siebes, Hayo Reijers and Hans Bodlaender. I would henceforth like to thank all of them for being part of my thesis committee.

Although this dissertation and the research therein was done by myself, this would never have been possible without the professional and personal support I have received from my supervisors, colleagues, friends and family.

First and foremost, I want to extend my gratitude to my supervisor, Wouter Swierstra, who was of central importance to this manuscript and, naturally, my PhD. Thank you for giving me the opportunity to work in such an interesting subject and for providing support whenever necessary, but at the same time giving me the space and freedom to grow and start learning how to be a better researcher. I also want to thank Gabriele Keller for accepting to be my *promotor* and for her invaluable help and advice with many bureaucratic aspects of being a PhD candidate. Additionally, I extend my gratitude to Pierre-Evariste Dagand for all the help he provided from technical issues to general good advice and the important feedback he gave on early versions of this manuscript, thank you for always asking the hard questions.

Next, I want to thank my (ex) lab-mates at BBG-5.70, Alejandro and João. They have welcomed me into Utrecht and into my PhD and have provided countless outstanding discussions, which were (mostly) scientifically relevant. Jurrian was probably our most common visitor at BBG-5.70 and he was always interested in what I was doing and did provide invaluable scientific advice with the addition of amazing band suggestions every now and then, thank you! I would also like to extend my gratitude to the entirety of the group and staff that routinely attended our Monday meetings – Johan, Wishnu, Jeroen, Raja, Nico, Sergey, Anna-Lena, Vedran, Trevor, Matthijs, Ivo, Iris, Isaac, Samira and Saba – for always being helpful, eager to discuss and listening to countless talk rehearsals on Mondays. Finally, I want to thank Arian and Martijn who were also frequent visitor at BBG-5.70 and always brought interesting topics to discuss and amazing Linux tricks. Finally, I would like to thank Lisa and Vedran for being my paranymphs and helping to organize the different aspects associated with the defense of this manuscript.

Subsequently, I would also like to thank my girlfriend, Nadine, for all her love, support and encouragement. Thank you for being there to help me get through all the frustrations and see the light at the end of the tunnel during the later years of my PhD. I can't wait to spend the years after my PhD with you by my side.

Next I want to thank Lisa for being instrumental to my happiness while in Utrecht, and consequently, to this thesis. Thank you for your support, companionship, for al-

ways being there for one last beer and for organizing so many great activities with me. I cherish all of our shared experiences very much and I am looking forward to the ones that are yet to come. If I may, let me propose Vietnamese distilleries for a strong next candidate!

I have been blessed with great friends here in Utrecht and they have all contributed to the completion of this thesis by being there. Thank you Vedran and Elena for so many kayaking and pandemic afternoons; Thank you Antonello for the yearly Eurovision gatherings and great cheeses; Thank you Kristhell and Ivan for the yearly Sinterklass; Thank you Marina for believing in my cooking career; Thank you Lisa for the brewery visits; Thank you João for introducing me to the ING; Thank you Joel, Seraina, Ivica, Laura, Jurica and all the people that have brought nothing but happiness, contributing greatly to my social well-being during my PhD.

I only met most of these amazing people because I went to Jan Primus one given evening, for the Social Wednesday event. Little did I expect I would be coming back pretty much every subsequent Wednesday and meet virtually all my friends there. Jan Primus certainly was our second home for a little while. Fast forward a few years and I had the privilege of being part of the activity committee of the ING alongside Lisa and João. Lisa, João and I then started hosting the Social Wednesdays and many other events ourselves. This would not have been possible if Laura and Justus were not behind the bar, welcoming us and making Jan Primus our second home in Utrecht, thank you. I would also like to thank the board of the International Neighbor Group for giving me this opportunity.

Last but not least, I want to thank my family in this last paragraph, were I'll switch to Portuguese. Obrigado Mãe, Pai e Irmão pelo amor e orgulho. Obrigado por sempre prover o ambiente e suporte necessário para que eu pudesse algum dia chegar onde eu cheguei. O único aspecto negativo desta experiência toda é que infelizmente vivemos muito longe. Com saudades e carinho, um beijo.



Utrecht Universiteit
Department of Computing Sciences