

Best Title in the Universe

42

Victor Cacciari Miraldo Wouter Swierstra

University of Utrecht

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

Abstract

stuff

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

General Terms Haskell

Keywords Haskell

1. Introduction

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

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

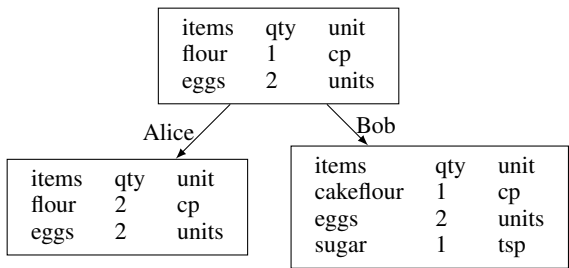


Figure 1. Sample CSV files

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

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

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

Figure 2. Output from `diff3`

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

We propose a structural diff that is not only generic but also able to track changes in a way that the user has the freedom to decide which is the fundamental editing unit. Our work is built on top of [FIXBIB], but we extend it in order to handle merging of patches. We also propose extensions to this algorithm capable of detecting purely structural operations such as swapping and cloning.

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

Contributions

Background

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

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, or to redistribute to lists, contact the Owner/Author. Request permissions from permissions@acm.org or Publications Dept., ACM, Inc., fax +1 (212) 869-0481. Copyright held by Owner/Author. Publication Rights Licensed to ACM.

2. Structural Diffing

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

- Give some context: *Tree-edit distance*;

To Research!

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

2.1 Context Free Datatypes

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

```
data U : N → Set where
  u0  : {n : N} → U n
  u1  : {n : N} → U n
  _⊕_ : {n : N} → U n → U n → U n
  _⊗_ : {n : N} → U n → U n → U n
  β   : {n : N} → U (suc n) → U n → U n
  μ   : {n : N} → U (suc n) → U n
  v1  : {n : N} → U (suc n)
  wk  : {n : N} → U n → U (suc n)
```

- Explain the patching problem.
- We want a type-safe approach.
- Argue that the types resulting from our parser are in a sub-language of what we treated next.

Having a regular, yet extensible, way to parse different files gives us a stepping stone to start discussing how to track differences in the results of those parsers. There has been plenty of research on this topic ([CITE STUFF!]), however, most of the time data is converted to an untyped intermediate representation. We would like to stay type-safe. [WHY?] Our research shows that the type $D\ a$ that expresses the differences between two elements of type a can be determined by induction on a 's structure.

- introduce edit-script, diffing and patching or apply

For example, let us see the differences between the original CSV and Alice's edits:

items	,qty	,unit
wheat-flour	,2	,cp
eggs	,2	,units
Alice		
items	,qty	,unit
wheat-flour	,1	,cp
eggs	,2	,units
Original		

Figure 3. Alice's edits

As we can see, Alice edited the	Cpy "items,qty,unit" (Del "weat-flour,1,cp" (Ins "wheat-flour,2,cp" (Cpy "eggs,2,unis" End)))
---------------------------------	---

Figure 4. Line-based edit-script for figure 3

the second line of the file. A line-based edit script, which is something that transform a file into another, would look like the one presented in figure 4. That edit script has a few problems: (A) it is deleting and inserting almost identical lines and (B) it is unaware of the CSV file semantics, making it harder to identify actual conflicts.

3. Sharing of Recursive Subterms

- If we want to be able to share recursive subexpressions we need a mutually recursive approach.

4. Remarks on Type Safety

- At which level of our design space we would like type-safety?
- Maybe after introducing the matrix idea it is clear that type-safety might be desirable only on the diff level, not on the patch level.

5. Related Work

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

6. Conclusion

- This is what we take out of it.