Structured diffs: theory and practice ICFP PC @ SLC

Victor Cacciari Miraldo, Pierre-Evariste Dagand, Giovanni Garufi, Marco Vassena and Wouter Swierstra



Universiteit Utrecht

The diff utility

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

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



Example: comparing two files

slc-teams.csv

Real Salt Lake, Soccer

Utah Jazz, Basketball

Salt Lake Bees, Baseball



Example: comparing two files

slc-teams.csv

Real Salt Lake, Soccer

Utah Jazz, Basketball

Salt Lake Bees, Baseball

slc-teams-fixed.csv

Real Salt Lake, Football

Utah Jazz, Basketball

Salt Lake Bees, Baseball



Example: comparing two files

-Real Salt Lake, Soccer +Real Salt Lake, Football Utah Jazz, Basketball Salt Lake Bees, Baseball

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



Smallest edit script

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



Universiteit Utrecht

Smallest edit script

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

In other words, it tries to preserve as much information as possible.



Smallest edit script

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

In other words, it tries to preserve as much information as possible.

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



slc-teams-fixed.csv

Real Salt Lake, Football Utah Jazz, Basketball Salt Lake Bees, Baseball

How would this file change if I add a new column?



Universiteit Utrecht

-Real Salt Lake, Football
+Real Salt Lake, Football, 2004
-Utah Jazz, Basketball
+Utah Jazz, Basketball, 1979
-Salt Lake Bees, Baseball
+Salt Lake Bees, Baseball, 1994



-Real Salt Lake, Football
+Real Salt Lake, Football, 2004
-Utah Jazz, Basketball
+Utah Jazz, Basketball, 1979
-Salt Lake Bees, Baseball
+Salt Lake Bees, Baseball, 1994

Adding a new column changes **every line** in our original file. Where conceptually, we are not modifying any existing **data**.



Universiteit Utrecht

-Real Salt Lake, Football
+Real Salt Lake, Football, 2004
-Utah Jazz, Basketball
+Utah Jazz, Basketball, 1979
-Salt Lake Bees, Baseball
+Salt Lake Bees, Baseball, 1994

Adding a new column changes **every line** in our original file. Where conceptually, we are not modifying any existing **data**. Not all data is best represented by a list of lines!

This is particularly important when using diff to compare *source code*.

What is the diff over structured data?



Universiteit Utrecht

Questions

- How can we represent a family of data types?
- How can we represent patches on these data types?
- Does this give a better account of software evolution?



Questions

How can we represent a family of data types?

- How can we represent patches on these data types?
- Does this give a better account of software evolution?



Universiteit Utrecht

Universe of discourse

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

```
data Atom : Set where
```

- K : U -> Atom
- I : Atom

```
Prod : <mark>Set</mark>
Prod = List Atom
```

```
Sum : Set
Sum = List Prod
```

Here we assume some 'base universe' **U**, storing the atomic types such as integers, characters, etc.



Semantics

We can interpret these types as *pattern functors*:

```
elA : Atom -> (Set -> Set)
elA I X = X
elA (K u) X = elU u
elP : Prod -> (Set -> Set)
elP[] X = Unit
elP (a :: as) X = Pair (elA alpha X) (elP pi X)
elS : Sum -> (Set -> Set)
elS[] X = Empty
elS (p :: ps) X = Either (elP p X) (elS ps X)
```



Universiteit Utrecht

Fixpoints

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

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

data Fix (s : Sum) : Set where
 <_> : elS s (Fix s) -> Fix s



Example: 2-3 trees

We can represent 2-3-trees defined as follows:

data Tree	:	Set	whe	ere							
leaf	:	Tree									
2-node	:	Nat	->	Tree	->	Tree	->	Tree			
3-node	:	Nat	->	Tree	->	Tree	->	Tree	->	Tree	

by the following sum-of-products:



Universiteit Utrecht

Questions

How can we represent a family of data types?

How can we represent patches on these data types?

Does this give a better account of software evolution?



2-3-trees

treeA = 2-node 7 t1 t2

treeB = 3-node 12 (2-node 7 t1 leaf) leaf leaf

What edit script should transform treeA to treeB?



Universiteit Utrecht

2-3-trees

treeA = 2-node 7 t1 t2

treeB = 3-node 12 (2-node 7 t1 leaf) leaf leaf

What edit script should transform treeA to treeB?

It is not just a list of insertions and deletions!

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

We will use a type indexed data type to account for changes.



Representing diffs

Our universe consists of three separate layers:

- sums
- products
- atomic values

We'll define what it means to modify each of these layers – from these pieces we can define our overall type for diffs.



Universiteit Utrecht

Spines: changes to sums

Given two arbitrary tree structures, **x** and **y**, we can identify the following three cases:

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



Spines: changes to sums

Given two arbitrary tree structures, **x** and **y**, we can identify the following three cases:

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

To represent patches, we need a data type that describes these three cases.

But what information should each constructor record?



Spines

Assuming that we know what patches on atoms (pAt) and products (pAl) are we can define:

```
data S (\sigma : Sum) : Set where

Scp : S \sigma

Scns : (C : Constr \sigma)

-> All pAt (fields C)

-> S \sigma

Schg : (C1 C2 : Constr \sigma)

-> pAl (fields C1) (fields C2)

-> S \sigma
```

We still need to define how to diff products and atoms.



If we have reconciled the choice of constructor, how to we compare the constructor fields?



Universiteit Utrecht

If we have reconciled the choice of constructor, how to we compare the constructor fields?

Each value constructed in our universe has a *list of fields* – the product structure.

Given two such lists, we need to compare them somehow.

Yet these fields may store values of very different types!



If we have reconciled the choice of constructor, how to we compare the constructor fields?

Each value constructed in our universe has a *list of fields* – the product structure.

Given two such lists, we need to compare them somehow.

Yet these fields may store values of very different types!

The good news, however, is that we can reuse ideas from the classic **diff** algorithm at this point.



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

- copies over fields;
- deletes fields;
- inserts new fields.



Alignments

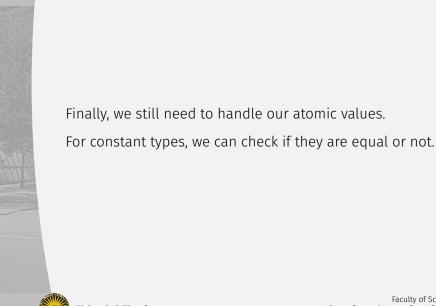
data Al : Prod \rightarrow Prod \rightarrow Set where A0 : Al At [] [] AX : At $\alpha \rightarrow$ Al $\pi 2 \ \pi 1 \rightarrow$ Al At ($\alpha :: \pi 2$) ($\alpha :: \pi 1$) Adel : elA a \rightarrow Al $\pi 2 \ \pi 1 \rightarrow$ Al ($\alpha :: \pi 2$) $\pi 1$ Ains : elA a \rightarrow Al $\pi 2 \ \pi 1 \rightarrow$ Al $\pi 2 \ (\alpha :: \pi 1)$

A value of type Al $\pi 2 \pi 1$ prescribes which fields of one constructor are matched with which fields of another.



Universiteit Utrecht

Atoms





Atoms



Universiteit Utrecht

Finally, we still need to handle our atomic values.

But what about recursive subtrees?

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

Faculty of Science Information and Computing Sciences

23

Handling recursive data types

So far our spines compare the outermost constructors.

Oftentimes, you may want to delete certain constructors (exposing its subtrees) or insert new constructors.

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



Accounting for recursion

Our final patch type identifies three cases:

- 1. The insertion of a new constructor, together with all-but-one of its fields;
- 2. The deletion of the outermost constructor, together with all-but-one of its fields;
- 3. A choice of spine, alignment, and a patch on atomic values;

The first two require additional information – a context – to point out *where* to insert/delete a subtree.



Applying patches

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

```
apply : Patch \rightarrow Fix \sigma \rightarrow Maybe (Fix \sigma)
```

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

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



Questions

- How can we represent a family of data types?
- How can we represent patches on these data types?
- Does this give a better account of software evolution?



Case study: Clojure

- We've instantiated this algorithm to a simplified Clojure AST in Haskell;
- By implementing a simple Clojure parser, we can now compare Clojure programs.



Case study: Clojure

- We've instantiated this algorithm to a simplified Clojure AST in Haskell;
- By implementing a simple Clojure parser, we can now compare Clojure programs.
- And by mining the commit history of the top Clojure repositories on GitHub, we can try to quantify the performance our algorithm.



Universiteit Utrecht



Collect data

- A) False conflicts two changes to the same line that do not overlap in the AST
- B) Fixable conflicts two changes to the same atomic value, where knowing the abstract syntax tree allows us to resolve them automatically/interactively.
- C) *True conflicts* two atomic values (integers, variables, etc.) changed in different ways



Contributors	LOC	Commits	Conflicts	А	В	С
35	14,693	2,416	18	8	2	8
62	4,557	1,064	17	6	5	6
66	9,370	1,271	8	2	2	4
43	8,028	1,366	46	17	14	15
109	5,193	1,111	5	1	0	4
36	6,604	1,818	12	1	4	7
65	28,790	586	12	3	2	7
33	803	227	1	1	0	0
92	894	18,857	27	5	2	20
82	16,478	1,282	40	9	22	9
47	1,099	401	4	2	0	2
86	6,515	1,464	6	4	0	2
315	10,669	4,484	28	12	4	12
42	2,965	347	8	6	1	1
46	23,778	6,641	90	46	11	33
55	27,935	2,996	50	21	6	23
59	1,1206	1,403	24	13	5	6
34	1,341	960	10	1	8	3
114	16,586	1,654	6	3	0	3
99	4,909	958	40	4	31	5
			452	165	117	170
	35 62 66 43 109 36 65 33 92 82 47 86 315 42 46 55 59 34 114	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$

Results



Interpreting these results

 Conflicts are rare! 452 conflicts found in tens of thousands of commits.



Universiteit Utrecht

Interpreting these results

 Conflicts are rare! 452 conflicts found in tens of thousands of commits.

Or perhaps hard to observe as rebasing can rewrite history, complicated pull requests abandoned, etc.

Structure aware algorithms can beat line-based diff



Interpreting these results

 Conflicts are rare! 452 conflicts found in tens of thousands of commits.

Or perhaps hard to observe as rebasing can rewrite history, complicated pull requests abandoned, etc.

Structure aware algorithms can beat line-based diff

But performance of our algorithm is still lagging behind.



Questions?



Universiteit Utrecht