

# PROOF REPAIR

TALIA RINGER

A dissertation  
submitted in partial fulfillment of the  
requirements for the degree of

Doctor of Philosophy

University of Washington

2021

Reading Committee:

TODO, Chair

TODO

TODO

Program Authorized to Offer Degree:  
Computer Science & Engineering



© Copyright 2021

Talia Ringer



ABSTRACT

PROOF REPAIR

Talia Ringer

Chairs of the Supervisory Committee:

TODO

Computer Science & Engineering

Abstract will go here.



To my family.



I love all of you.





---

## CONTENTS

---

1	INTRODUCTION	3
2	MOTIVATING PROOF REPAIR	5
2.1	Proof Development	5
2.2	Proof Maintenance	5
2.3	Proof Repair	6
3	PROOF REPAIR BY EXAMPLE	7
3.1	Motivating Example	7
3.2	Approach	9
3.3	Differencing	12
3.4	Transformation	13
3.5	Implementation	14
3.6	Results	18
3.7	Conclusion	24
4	PROOF REPAIR ACROSS TYPE EQUIVALENCES	25
4.1	Motivating Example	25
4.2	Approach	27
4.3	Differencing	30
4.4	Transformation	37
4.5	Implementation	37
4.6	Results	42
4.7	Conclusion	47
5	RELATED WORK	49
5.1	Programs	49
5.2	Proofs	50
6	CONCLUSIONS & FUTURE WORK	55



---

## ACKNOWLEDGMENTS

---

I’ve always believed the acknowledgments section to be one of the most important parts of a paper. But there’s never enough room to thank everyone I want to thank. Now that I have the chance—where do I begin?

We got other wonderful feedback on the paper from Cyril Cohen, Tej Chajed, Ben Delaware, Jacob Van Geffen, Janno, James Wilcox, Chandrakana Nandi, Martin Kellogg, Audrey Seo, James Decker, and Ben Kushigian. And we got wonderful feedback on e-graph integration for future work from Max Willsey, Chandrakana Nandi, Remy Wang, Zach Tatlock, Bas Spitters, Steven Lyubomirsky, Andrew Liu, Mike He, Ben Kushigian, Gus Smith, and Bill Zorn. The Coq developers have for years given us frequent and efficient feedback on plugin APIs for tool implementation. Merge in more PUMPKIN Pi thank yous. Merge in survey paper thank yous.

Dan Grossman, Jeff Foster, Zach Tatlock, Derek Dreyer, Alexandra Silva, the Coq community (Emilio J. Gallego Arias, Enrico Tassi, Gaëtan Gilbert, Maxime Dénès, Matthieu Sozeau, Vincent Laporte, Théo Zimmermann, Jason Gross, Nicolas Tabareau, Cyril Cohen, Pierre-Marie Pédro, Yves Bertot, Tej Chajed, Ben Delaware, Janno), coauthors, Valentin Robert, my family, PLSE lab (especially Chandrakana Nandi oh my gosh), James Wilcox, Jasper Hugunin, Marisa Kirisame, Jacob Van Geffen, Martin Kellogg, Audrey Seo, James Decker, Ben Kushigian, Gus Smith, Max Willsey, Zach Tatlock, Steven Lyubomirsky, Andrew Liu, Mike He, Ben Kushigian, Bill Zorn, Anders Mörtberg, Conor McBride, Carlo Angiuli, Bas Spitters, UCSD Programming Systems group, Misha, PL Twitter, Roy, Vikram, Alex Polozov, Esther, Ellie, Mer, students, Qi, Saba.



---

## INTRODUCTION

---

Motivation for verifying systems

Era of scale—enter proof engineering [70]

Looking back (Social Processes [30]), development has come a long way, but maintenance is still hard! And this is a problem in practice!

But missed opportunity: automation doesn't understand that proofs evolve

So we build automation that does, and we call this proof repair. Proof repair shows that there is reason to believe that verifying a modified system should often, in practical use cases, be easier than verifying the original the first time around.

Or, in other words (thesis statement): Changes in programs, specifications, and proofs carry information that a tool can extract, generalize, and apply to fix other proofs broken by the same change. A tool that automates this can save work for proof engineers relative to reference manual repairs in practical use cases.

Key technical bit: differencing and program transformations, taking advantage of the rich and structured language proofs are written in.

We implement this in a tool suite for Coq, get some sweet results.

Pave path to the next era of verification

### READING GUIDE

How to read this thesis

Mapping of papers to chapters

Authorship statements for included paper materials, to credit coauthors

Expected reader background & where to find more info



# 2

---

## MOTIVATING PROOF REPAIR

---

Before we talk more about proof repair, it helps to know what it's like to develop and maintain proofs to begin with, and what happens under the hood when you do that. This chapter gives you that context, then explains the high-level approach to proof repair that builds on that.

### 2.1 PROOF DEVELOPMENT

Cartoon version of development: program, spec, proof

Proof assistants: short overview of foundations & different options (survey paper), then say focus on Coq

Slightly less brief overview of Coq and its foundations and automation and so on (including proof terms), going through a running example of proof development in Coq

### 2.2 PROOF MAINTENANCE

Problem is when something changes—change something in running example

There are a lot of development processes people use to make proofs less likely to break to begin with (survey paper)

But still, even with these, the reality: This happens all the time (REPLICA)

And in fact not just after developing a proof, but during development too (REPLICA)

And breaks proofs even for experts (REPLICA)

And it's an extra big problem when you have a large development and the changes are outside of your control

Hence Social Processes

Why automation breaks, even with good development processes

Hence proof repair—smarter automation

### 2.3 PROOF REPAIR

Name inspired by program repair, but quite different as we'll soon see.

Recall thesis: Changes in programs, specifications, and proofs carry information that a tool can extract, generalize, and apply to fix other proofs broken by the same change. A tool that automates this can save work for proof engineers relative to reference manual repairs in practical use cases.

Proof repair accomplishes this using a combination of differencing and program transformations.

Differencing extracts the information from the change in program, specification, or proof.

The transformations then generalize that information to a more general fix for other proofs broken by the same change.

The details of applying the fix vary by the kind of fix, as we'll soon see.

Crucially, all of this happens over the proof terms in this rich language we saw in the Development section. This is kind of the key insight that makes it all work.

This is great because this language gives us so much information and certainty. This helps us with two of the biggest challenges from program repair. (generals related work)

But it's also challenging because this language is so unforgiving. Plus, in the end, we need these tactic proofs, not just proof terms. So we can't just reuse program repair tools. (generals related work)

So next two chapters will show two tools in our tool suite that work this way, how they handle these challenges, and how they save work.



# 3

---

## PROOF REPAIR BY EXAMPLE

---

The first tool (PUMPKIN PATCH) focuses on changes in programs and specifications, though these changes are limited in scope as we'll see later.

What this tool does is, when programs and specifications change and this breaks a lot of proofs, it lets the proof engineer fix just one of those proofs. It then generalizes the example patch into something that can fix other proofs broken by the same change.

So in other words, the information from those changes is carried in the difference between the old and new version of the example patched proof. PUMPKIN PATCH generalizes that information.

Application can be automated in some cases at the end, or it can be manual.

The work saved is shown retroactively on case studies replaying changes from large proof developments in Git. Results for this tool are preliminary compared to what we'll see later, since this was the first prototype.

### 3.1 MOTIVATING EXAMPLE

Traditional proof automation considers only the current state of theorems, proofs, and definitions. This is a missed opportunity: verification projects are rarely static. Like other software, these projects evolve over time.

With traditional proof automation, the burden of change largely falls on proof engineers. This does not have to be true. Proof automation can view theorems, proofs, and definitions as fluid entities: when a proof or specification changes, a tool can search the difference between the old and new versions for a *reusable patch* that can fix broken proofs.

**WITHOUT PROOF REPAIR** Experienced Coq programmers use design principles and custom tactics to make proofs resilient to change. These techniques are useful for large proof developments, but they place the burden of change on the programmer. This can be problematic when change occurs outside of the programmer's control.

<pre> Definition IZR (z:Z) : R :=   match z with     Z0 =&gt; 0     Zpos n =&gt; INR (Pos.to_nat n)     Zneg n =&gt; - INR (Pos.to_nat n) end. </pre>	<pre> Definition IZR (z:Z) : R :=   match z with     Z0 =&gt; 0     Zpos n =&gt; IPR n     Zneg n =&gt; - IPR n end. </pre>
---	---

Figure 1: Old (left) and new (right) definitions of IZR in Coq. The old definition applies injection from naturals to reals and conversion of positives to naturals; the new definition applies injection from positives to reals.

Consider a commit from the Coq 8.7 release [59]. This commit redefined injection from integers to reals (Figure 1). This change broke 18 proofs in the standard library.

The Coq developer who committed the change fixed the broken proofs, then made an additional 12 commits to address the change in `coq-contribs`, a regression suite of projects that the Coq developers maintain as versions change. Many of these changes were simple. For example, the developer wrote a lemma that describes the change:

**Lemma** INR\_IPR :  $\forall p$ , INR (Pos.to\_nat p) = IPR p.

The developer then used this lemma to fix broken proofs within the standard library. For example, one proof broke on this line:

```
rewrite Pos2Nat.inj_sub by trivial.✗
```

It succeeded with the lemma:

```
rewrite <- 3!INR_IPR, Pos2Nat.inj_sub by trivial.✓
```

These changes are outside-facing: Coq users have to make similar changes to their own proofs when they update from Coq 8.6 to Coq 8.7. The Coq developer can update some tactics to account for this, but it is impossible to account for every tactic that users could use. Furthermore, while the developer responsible for the changes knows about the lemma that describes the change, the Coq user does not. The Coq user must determine how the definition has changed and how to address the change, perhaps by reading documentation or by talking to the developers.

**WITH PROOF REPAIR** When a user updates the Coq standard library, a proof repair tool can determine that the definition has changed, then analyze changes in the standard library and in `coq-contribs` that resulted from the change in definition (in this case, rewriting by the lemma). It can extract a reusable patch from those changes, which it can automatically apply within broken user proofs. The user never has to consider how the definition has changed.

### 3.2 APPROACH

In the example from Section 3.1, we can see how the example change in one proof carries enough information to fix other proofs broken by the same change (namely the rewrite by `INR\_IPR`). So a tool can extract that, generalize it, and use it to fix other proofs broken by the same change.

The key insight behind PUMPKIN’s approach is that this is true more generally. To use PUMPKIN, the programmer modifies a single proof script to provide an *example* of how to adapt a proof to a change. PUMPKIN extracts that information into a *patch candidate*—which is localized to the context of the example, but not enough to fix other proofs broken by the change. It then generalizes that candidate into a *reusable patch*: a function that can be used to fix other broken proofs broken by the same change, which PUMPKIN defines as a Coq term.

In other words, looking back to the thesis statement, the information shows up in the difference between versions of the example patched proof. PUMPKIN can extract and generalize that information. Application works with hint databases or is manual. Here is the system diagram for PUMPKIN. The PUMPKIN repository contains a detailed user guide.

As mentioned earlier, PUMPKIN does this using a combination of semantic differencing and program transformations. Differencing looks at the difference between versions of the example patched proof for this information, and finds the candidate. Then, program transformations modify that candidate to produce the reusable proof patch.

And of course all of this happens over proof terms, since tactics might hide necessary information. Of course this is hard to see on the example from Section 3.1, since we were lucky enough to see the difference in tactics here. Let’s look at a toy example for which that isn’t true.

To motivate this workflow, consider using PUMPKIN to search the proofs in Figure 2 for a patch between conclusions. Except we will show a place where the lemma is actually applied. Note that the tactics don’t change even though the terms do—and even though the change could break other proofs.

So what do we do? We invoke the plugin using `old` and `new` as the example change:

```
Patch Proof old new as patch.
```

PUMPKIN first determines the type that a patch from `new` to `old` should have. To determine this, it semantically *diffs* the types and finds this goal type (line 2):

```
∀ n m p, n <= m -> m <= p -> n <= p -> n <= p + 1
```

It then breaks each inductive proof into cases and determines an intermediate goal type for the candidate. In the base case, for example,

<pre> 1 Theorem old: ∀ (n m p : nat),   n &lt;= m -&gt; m &lt;= p -&gt; 2   n &lt;= p + 1.   (* P p *) 3 Proof. 4   intros. induction H0. 5   - auto with arith. 6   - constructor. auto. 7 Qed. 8 9 fun (n m p : nat) (H : n &lt;= m 10   ) (H0 : m &lt;= p) =&gt; 11   le_ind 12     m 13     (* 14       m *) 15     (fun p0 =&gt; n &lt;= p0 + 1) 16     (* P *) 17     (le_plus_trans n m 1 H) 18     (* : P m *) 19     (fun (m0 : nat) (_ : m &lt;= 20       m0) (IHle : n &lt;= m0 + 1) =&gt; 21         le_S n (m0 + 1) IHle) 22     p 23     (* 24       p *) 25     H0 </pre>	<pre> 1 Theorem new: ∀ (n m p : nat 2   ), n &lt;= m -&gt; m &lt;= p -&gt; 3   n &lt;= p.   (* P' p *) 4 Proof. 5   intros. induction H0. 6   - auto with arith. 7   - constructor. auto. 8 Qed. 9 10 fun (n m p : nat) (H : n &lt;= 11   m) (H0 : m &lt;= p) =&gt; 12   le_ind 13     m 14     (* m *) 15     (fun p0 =&gt; n &lt;= p0) 16     (* P' *) 17     H 18     (* : P' m *) 19     (fun (m0 : nat) (_ : m 20       &lt;= m0) (IHle : n &lt;= m0) =&gt; 21         le_S n m0 IHle) 22     p 23     (* p *) 24     H0 </pre>
--	--

Figure 2: Two proofs with different conclusions (top) and the corresponding proof terms (bottom) with relevant type information. We highlight the change in theorem conclusion and the difference in terms that corresponds to a patch.

it *diffs* the types and determines that a candidate between the base cases of *new* and *old* should have this type (lines 11 and 12):

```
(fun p0 => n <= p0) m -> (fun p0 => n <= p0 + 1) m
```

It then *diffs* the terms (line 13) for such a candidate:

```
fun n m p H0 H1 =>
  (fun (H : n <= m) => le_plus_trans n m 1 H)
: ∀ n m p, n <= m -> m <= p -> n <= m -> n <= m + 1
```

This candidate is close, but it is not yet a patch. This candidate maps base case to base case (it is applied to *m*); the patch should map conclusion to conclusion (it should be applied to *p*).

This is where the transformations come in. There are four:

1. *Patch specialization* to arguments
2. *Patch abstraction* of arguments or functions
3. *Patch inversion* to reverse a patch
4. *Lemma factoring* to break a term into parts

Here, PUMPKIN *abstracts* this candidate by *m* (line 11), which lifts it out of the base case:

```
fun n0 n m p H0 H1 =>
  (fun (H : n <= n0) => le_plus_trans n n0 1 H)
: ∀ n0 n m p, n <= m -> m <= p -> n <= n0 -> n <= n0 + 1
```

PUMPKIN then *specializes* this candidate to *p* (line 16), the argument to the conclusion of *le\_ind*. This produces a patch:

```
patch n m p H0 H1 :=
  (fun (H : n <= p) => le_plus_trans n p 1 H)
: ∀ n m p, n <= m -> m <= p -> n <= p -> n <= p + 1
```

The user can then use *patch* to fix other broken proofs. For example, given a proof that applies *old*, the user can use *patch* to prove the same conclusion by applying *new*:

```
apply old.✓
apply patch. apply new.✓
```

This can happen automatically through hint databases.

This simple example uses only two transformations. The other transformations help turn candidates into patches in similar ways. We discuss all of this in detail later.

**CONFIGURATION** The components come together to form a proof patch finding procedure:

---

**Pseudocode:** find\_patch(term, term', direction)

---

- 1: *diff* types of term and term' for goals
  - 2: *diff* term and term' for candidates
  - 3: **if** there are candidates **then**
  - 4:   *factor, abstract, specialize, and/or invert* candidates
  - 5:   **if** there are patches **then return** patches
  - 6: **return** failure
-

PUMPKIN infers a *configuration* from the example change. This configuration customizes the highlighted lines for an entire class of changes: It determines what to diff on lines 1 and 2, and how to use the components on line 4.

For example, to find a patch for Figure 2, PUMPKIN used the configuration for changes in conclusions of two proofs that induct over the same hypothesis. Given two such proofs:

$$\begin{array}{l} \forall x, H\ x \rightarrow P\ x \\ \forall x, H\ x \rightarrow P'\ x \end{array}$$

PUMPKIN searches for a patch with this type:

$$\forall x, H\ x \rightarrow P'\ x \rightarrow P\ x$$

using this configuration:

---

```
1: diff conclusion types for goals
2: diff conclusion terms for candidates
3: if there are candidates then
4:   abstract and then specialize candidates
```

---

Later we will see real-world examples that demonstrate more configurations.

### 3.3 DIFFERENCING

The tool should be able to identify the semantic difference between terms. The semantic difference is the difference between two terms that corresponds to the difference between their types. Consider the base case terms in Figure 2 (line 13):

$$\begin{array}{l} \text{le\_plus\_trans } n\ m\ 1\ H : n \leq m + 1 \\ \text{le\_plus\_trans } n\ m\ 1\ H : n \leq m \end{array}$$

The semantic differencing component first identifies the difference in their types, or the *goal type*:

$$n \leq m \rightarrow n \leq m + 1$$

It then finds a difference in terms that has that type:

$$\text{fun } (H : n \leq m) \Rightarrow \text{le\_plus\_trans } n\ m\ 1\ H$$

This is the *candidate* for a reusable patch that the other components modify to find a patch.

Differencing operates over terms and types. Differencing tactics is insufficient, since tactics and hints may mask patches (line 5).<sup>1</sup> Furthermore, differencing is aware of the semantics of terms and types. Simply exploring the syntactic difference makes it hard to identify which changes are meaningful. For example, in the inductive case (line 14), the inductive hypothesis changes:

$$\begin{array}{l} \dots (\text{IHle} : n \leq m0 + 1) \dots \\ \dots (\text{IHle} : n \leq m0) \dots \end{array}$$

<sup>1</sup> Since this is a simple example, replaying an existing tactic happens to work. There are additional examples in the repository (Cex.v).

However, the type of `IH1e` changes for *any* two inductive proofs over `1e` with different conclusions. A syntactic differencing component may identify this change as a candidate. Our semantic differencing component knows that it can ignore this change.

Plus parts of Inside the Core, Testing Boundaries, Future Work

How differencing works in detail

Limitations and whether they're addressed in later tools yet or not

### 3.4 TRANSFORMATION

**PATCH SPECIALIZATION** The tool should be able to specialize a patch candidate to specific arguments as determined by the differences in terms. To find a patch for Figure 2, for example, PUMPKIN must specialize the patch candidate to `p` to produce the final patch.

**PATCH ABSTRACTION** A tool should be able to abstract patch candidates of this form by the common argument:

```
candidate : P' t -> P t
candidate_abs : ∀ t0, P' t0 -> P t0
```

and it should be able to abstract patch candidates of this form by the common function:

```
candidate : P t' -> P t
candidate_abs : ∀ P0, P0 t' -> P0 t
```

This is necessary because the tool may find candidates in an applied form. For example, when searching for a patch between the proofs in Figure 2, PUMPKIN finds a candidate in the difference of base cases. To produce a patch, PUMPKIN must abstract the candidate by the argument `m`. Abstracting candidates is not always possible; abstraction will necessarily be a collection of heuristics.

**PATCH INVERSION** The tool should be able to invert a patch candidate. This is necessary to search for isomorphisms. It is also necessary to search for implications between propositionally equal types, since candidates may appear in the wrong direction. For example, consider two list lemmas (we write `length` as `len`):

```
old : ∀ l' l, len (l' ++ l) = len l' + len l
new : ∀ l' l, len (l' ++ l) = len l' + len (rev l)
```

If PUMPKIN searches the difference in proofs of these lemmas for a patch from the conclusion of `new` to the conclusion of `old`, it may find a candidate *backwards*:

```
candidate l' l (H : old l' l) :=
  eq_ind_r ... (rev_length l)
: ∀ l' l, old l' l -> new l' l
```

The component can invert this to get the patch:

```
patch l' l (H : new l' l) :=
  eq_ind_r ... (eq_sym (rev_length l))
: ∀ l' l, new l' l -> old l' l
```

We can then use this patch to port proofs. For example, if we add this patch to a hint database [1], we can port this proof:

```
Theorem app_rev_len : ∀ l l',
  len (rev (l' ++ l)) = len (rev l) + len (rev l').
Proof.
  intros. rewrite rev_app_distr. apply old. ✓
Qed.
```

to this proof:

```
Theorem app_rev_len : ∀ l l',
  len (rev (l' ++ l)) = len (rev l) + len (rev l').
Proof.
  intros. rewrite rev_app_distr. apply new. ✓
Qed.
```

Rewrites like `candidate` are *invertible*: We can invert any rewrite in one direction by rewriting in the opposite direction. In contrast, it is not possible to invert the patch PUMPKIN found for Figure 2. Inversion will necessarily sometimes fail, since not all terms are invertible.

**LEMMA FACTORING** The tool should be able to factor a term into a sequence of lemmas. This can help break other problems, like abstraction, into smaller subproblems. It is also necessary to invert certain terms. Consider inverting an arbitrary sequence of two rewrites:

$$t := \text{eq\_ind\_r } G \dots (\text{eq\_ind\_r } F \dots)$$

We can view  $t$  as a term that composes two functions:

$$\begin{aligned} g &:= \text{eq\_ind\_r } G \dots \\ f &:= \text{eq\_ind\_r } F \dots \\ t &:= g \circ f \end{aligned}$$

The inverse of  $t$  is the following:

$$t^{-1} := f^{-1} \circ g^{-1}$$

To invert  $t$ , PUMPKIN identifies the factors  $[f; g]$ , inverts each factor to  $[f^{-1}; g^{-1}]$ , then folds and applies the inverse factors in the opposite direction.

plus parts of PUMPKIN PATCH Inside the Core, Testing Boundaries, Future Work

How the four transformations work in detail

Limitations and whether they're addressed in later tools yet or not

### 3.5 IMPLEMENTATION

parts of PUMPKIN PATCH Inside the Core, plus more

#### 3.5.1 Tool Details

While our system is a very early prototype under active development, we have made the source code available on Github.<sup>2</sup> The interested

<sup>2</sup> <http://github.com/uwplse/PUMPKIN-PATCH/tree/cpp18>



reader can follow along in the repository. Our prototype has no impact on the trusted computing base (Section 3.5.2.1).

### 3.5.1.1 Semantic Differencing

We implement semantic differencing over *trees*: PUMPKIN compiles each proof term into a tree (`evaluation.ml`). In these trees, every node is a type context, and every edge is an extension to that type context with a new term.<sup>3</sup> Correspondingly, type differencing (to identify goal types) compares nodes, and term differencing (to find candidates) compares edges.

The component (`differencing.ml`) uses these nodes and edges to prioritize semantically relevant differences. At the lowest level, it calls a primitive differencing function which checks if it can substitute one term within another term to find a function between their types.

The key benefit to this model is that it gives us a natural way to express inductive proofs, so that differencing can efficiently identify good candidates. Consider, for example, searching for a patch between conclusions of two inductive proofs of theorems about the natural numbers:

```
nat_ind P ... (fun (IH : P n) => ...) : ∀ n, P n
nat_ind P' ... (fun (IH : P' n) => ...) : ∀ n, P' n
```

In each case, the component diffs the terms in the dotted edges of the tree for `nat_ind` (Figure 3) to try to find a term that maps between conclusions of that case:

```
P' 0 -> P 0 (* base case candidate *)
P' (S n) -> P (S n) (* inductive case candidate *)
```

The component also knows that the change in the type of `IH` is inconsequential (it occurs for any change in conclusion). Furthermore, it knows that `IH` cannot show up as a hypothesis in the patch, so it attempts to remove any occurrences of `IH` in any candidate.

When the component finds a candidate, it knows `P'` and `P` as well as the arguments `0` or `(S n)`. This makes it simple to query abstraction for the final patch:

```
∀ n, P' n -> P n
```

The differencing component is *lazy*: it compiles terms into trees one step at a time. It then *expands* each tree as needed to find candidates (`expansion.ml`). For example, consider searching two functions for a patch between conclusions:

```
fun (t : T) => b
fun (t' : T) => b'
```

Differencing introduces a single term of type `T` to a common environment, then expands and recursively diffs the bodies `b` and `b'` in that environment.

<sup>3</sup> These trees are inspired by categorical models of dependent type theory [39].

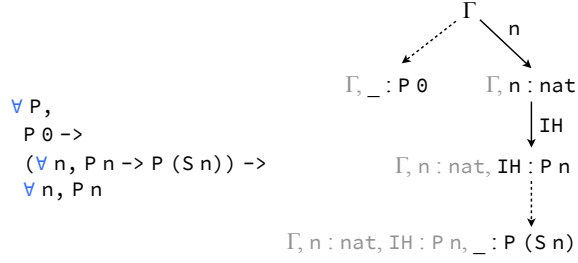


Figure 3: The type of (left) and tree for (right) the induction principle `nat_ind`. The solid edges represent hypotheses, and the dotted edges represent the proof obligations for each case in an inductive proof.

The tool always maintains pointers to easily switch between the tree and AST representations of the terms. This representation enables extensibility.

### 3.5.1.2 Transformations

**PATCH SPECIALIZATION** Specialization (`specialize.ml`) takes a patch candidate and some arguments, all of which are Coq terms. It applies the candidate to the arguments, then it  $\beta\iota$ -reduces [18] the result using Coq's `Reduction.nf_betaiota` function. It is the job of the patch finding procedure to provide both the candidate and the arguments.

**PATCH ABSTRACTION** Abstraction (`abstraction.ml`) takes a patch candidate, the goal type, and the function arguments or function to abstract. It first generalizes the candidate, wrapping it inside of a lambda from the type of the term to abstract. Then, it substitutes terms inside the body with the abstract term. It continues to do this until there is nothing left to abstract, then filters results by the goal type. Consider, for example, abstracting this candidate by `m`:

```
fun (H : n <= m) => le_plus_trans n m 1 H
: n <= m -> n <= m + 1
```

The generalization step wraps this in a lambda from some `nat`, the type of `m`:

```
fun (n0 : nat) =>
  (fun (H : n <= m) => le_plus_trans n m 1 H)
: ∀ n0, n <= m -> n <= m + 1
```

The substitution step replaces `m` with `n0`:

```
fun (n0 : nat) =>
  (fun (H : n <= n0) => le_plus_trans n n0 1 H)
: ∀ n0, n <= n0 -> n <= n0 + 1
```

Abstraction uses a list of *abstraction strategies* to determine what subterms to substitute. In this case, the simplest strategy works: The tool replaces all terms that are convertible to the concrete argument

$m$  with the abstract argument  $n0$ , which produces a single candidate. Type-checking this candidate confirms that it is a patch.

In some cases, the simplest strategy is not sufficient, even when it is possible to abstract the term. It may be possible to produce a patch only by abstracting *some* of the subterms convertible to the argument or function (we show an example of this in Section ??), or the term may not contain any subterms convertible to the argument or function at all. We implement several strategies to account for this. The combinations strategy, for example, tries all combinations of substituting only some of the convertible subterms with the abstract argument. The pattern-based strategy substitutes subterms that match a certain pattern with a term that corresponds to that pattern.

It is the job of the patch finding procedure to provide the candidate and the terms to abstract. In addition, each configuration includes a list of strategies. The configuration for changes in conclusions, for example, starts with the simplest strategy, and moves on to more complex strategies only if that strategy fails. This design makes abstraction simple to extend with new strategies and simple to call with different strategies for different classes of changes.

**PATCH INVERSION** Patch inversion (`inverting.ml`) exploits symmetry to try to reverse the conclusions of a candidate patch. It first factors the candidate using the factoring component, then calls the primitive inversion function on each factor, then finally folds the resulting list in reverse. The primitive inversion function exploits symmetry. For example, equality is symmetric, so the component can invert any application of `eq_ind` or `eq_ind_r` (any rewrite). Indeed, `eq_ind` and `eq_ind_r` are inverses, and are related by symmetry:

```
eq_ind_r A x P (H : P x) y (H0 : y = x) :=
  eq_ind x (fun y0 : A => P y0) H y (eq_sym H0)
```

If inversion does not recognize that the type is symmetric, it swaps subterms and type-checks the result to see if it is an inverse.

**LEMMA FACTORING** The lemma factoring component (`factoring.ml`) searches within a term for its factors. For example, if the term composes two functions, it returns both factors:

```
t : X -> Z (* term *)
[f : X -> Y; g : Y -> Z] (* factors *)
```

In this case, the component takes the composite term and  $x$  as arguments. It first searches as deep as possible for a term of type  $x \rightarrow y$  for some  $y$ . If it finds such a term, then it recursively searches for a term with type  $y \rightarrow z$ . It maintains all possible paths of factors along the way, and it discards any paths that cannot reach  $z$ .

The current implementation can handle paths with more than two factors, but it fails when  $y$  depends on  $x$ . Other components may benefit from dependent factoring; we leave this to future work.

### 3.5.1.3 Inside the Procedure

The implementation (`patcher.ml4`) of the procedure from Section ?? starts with a preprocessing step which compiles the proof terms to trees (like the tree in Figure 3). It then searches for candidates one step at a time, expanding the trees when necessary.

The PUMPKIN prototype exposes the patch finding procedure to users through the Coq command `Patch Proof` [Proof](#). PUMPKIN automatically infers which configuration to use for the procedure from the example change. For example, to find a patch for the case study in Section 3.6.1, we used this command:

```
Patch Proof Old.unsigned_range unsigned_range as patch.
```

PUMPKIN analyzed both versions of `unsigned_range` and determined that a constructor of the `int` type changed (Figure 4), so it initialized the configuration for changes in constructors.

Internally, PUMPKIN represents configurations as sets of options, which it passes to the procedure. The procedure uses these options to determine how to compose components (for example, whether to abstract candidates) and how to customize components (for example, whether semantic differencing should look for an intermediate lemma). To implement new configurations for different classes of changes, we simply tweak the options.

### 3.5.2 Workflow Integration

Needed: hints and so on, any work done since, the Git interface, whatever.

#### 3.5.2.1 Trusted Computing Base

A common concern for Coq plugins is an increase in the trusted computing base. The Coq developers provide a safe plugin API in Coq 8.7 to address this [31]. Our prototype takes this into consideration: While PUMPKIN does not yet support Coq 8.7, it only calls the internal Coq functions that the developers plan to expose in the safe API [47]. Furthermore, Coq type-checks terms that plugins produce. Since PUMPKIN does not modify the type checker, it cannot produce an ill-typed term.

## 3.6 RESULTS

Needed: key technical results

We used the PUMPKIN prototype to emulate three motivating scenarios from real-world code:

1. **Updating definitions** within a project  
(CompCert, Section 3.6.1)

<pre>Record int : Type :=   mkint { intval: Z; intrange         : 0 &lt;= intval &lt; modulus         }.</pre>	<pre>Record int : Type :=   mkint { intval: Z; intrange         : -1 &lt; intval &lt; modulus         }.</pre>
--	--

Figure 4: Old (left) and new (right) definitions of `int` in CompCert.

2. **Porting definitions** between libraries  
(Software Foundations, Section 3.6.2)
3. **Updating proof assistant versions**  
(Coq Standard Library, Section 3.6.3)

The code we chose for these scenarios demonstrated different classes of changes. For each case, we describe how PUMPKIN configures the procedure to use the core components for that class of changes. Our experiences with these scenarios suggest that patches are useful and that the components are effective and flexible.

**IDENTIFYING CHANGES** We identified Git commits from popular Coq projects that demonstrated each scenario. These commits updated proofs in response to breaking changes. We emulated each scenario as follows:

1. *Replay* an example proof update for PUMPKIN
2. *Search* the example for a patch using PUMPKIN
3. *Apply* the patch to fix a different broken proof

Our goal was to simulate incremental use of a patch finding tool, at the level of a small change or a commit that follows best practices. We favored commits with changes that we could isolate. When isolating examples for PUMPKIN, we replayed changes from the bottom up, as if we were making the changes ourselves. This means that we did not always make the same change as the user. For example, the real change from Section 3.6.1 updated multiple definitions; we updated only one.

PUMPKIN is a proof-of-concept and does not yet handle some kinds of proofs. In each scenario, we made minor modifications to proofs so that we could use PUMPKIN (for example, using induction instead of destruction). PUMPKIN does not yet handle structural changes like adding constructors or parameters, so we focused on changes that preserve structure, like modifying constructors. Chapter 4 describes an extension to PUMPKIN that supports changes in structure.

### 3.6.1 Updating Definitions

Coq programmers sometimes make changes to definitions that break proofs within the same project. To emulate this use case, we identified

<pre> Fixpoint bin_to_nat (b : bin) :   nat :=   match b with     B0 =&gt; 0     B2 b' =&gt; 2 * (bin_to_nat b     ')     B21 b' =&gt; 1 + 2 * (     bin_to_nat b')   end. </pre>	<pre> Fixpoint bin_to_nat (b : bin) :   nat :=   match b with     B0 =&gt; 0     B2 b' =&gt; (bin_to_nat b') +     (bin_to_nat b')     B21 b' =&gt; S ((bin_to_nat b     ') + (bin_to_nat b'))   end. </pre>
---	--

Figure 5: Definitions of `bin_to_nat` for Users A (left) and B (right).

a CompCert commit [50] with a breaking change to `int` (Figure 4). We used PUMPKIN to find a patch that corresponds to the change in `int`. The patch PUMPKIN found fixed broken inductive proofs.

**REPLAY** We used the proof of `unsigned_range` as the example for PUMPKIN. The proof failed with the new `int`:

```

Theorem unsigned_range:
  ∀(i : int), 0 <= unsigned i < modulus.
Proof.
  intros i. induction i using int_ind; auto.✗

```

We replayed the change to `unsigned_range`:

```

intros i. induction i using int_ind. simpl. omega.✓

```

**SEARCH** We used PUMPKIN to search the example for a patch that corresponds to the change in `int`. It found a patch with this type:

```

∀ z : Z, -1 < z < modulus -> 0 <= z < modulus

```

**APPLY** After changing the definition of `int`, the proof of the theorem `repr_unsigned` failed on the last tactic:

```

Theorem repr_unsigned:
  ∀(i : int), repr (unsigned i) = i.
Proof.
  ... apply Zmod_small; auto.✗

```

Manually trying `omega`—the tactic which helped us in the proof of `unsigned_range`—did not succeed. We added the patch that PUMPKIN found to a hint database. The proof of the theorem `repr_unsigned` then went through:

```

... apply Zmod_small; auto.✓

```

### 3.6.1.1 Configuration

This scenario used the configuration for changes in constructors of an inductive type. Given such a change:

```

Inductive T := ... | C : ... -> H -> T
Inductive T' := ... | C : ... -> H' -> T'

```

PUMPKIN searches two inductive proofs of theorems:

```

∀ (t : T), P t
∀ (t : T'), P t

```

for an isomorphism<sup>4</sup> between the constructors:

```
... -> H -> H'
... -> H' -> H
```

The user can apply these patches within the inductive case that corresponds to the constructor *C* to fix other broken proofs that induct over the changed type. PUMPKIN uses this configuration for changes in constructors:

- 
- 1: *diff* inductive constructors for goals
  - 2: use *all components* to recursively search for changes in conclusions of the corresponding case of the proof
  - 3: **if** there are candidates **then**
  - 4:   try to *invert* the patch to find an isomorphism
- 

### 3.6.2 Porting Definitions

Coq programmers sometimes port theorems and proofs to use definitions from different libraries. To simulate this, we used PUMPKIN to port two solutions [2, 7] to an exercise in Software Foundations to each use the other solution’s definition of the fixpoint `bin_to_nat` (Figure 5). We demonstrate one direction; the opposite was similar.

**REPLAY** We used the proof of `bin_to_nat_pres_incr` from User A as the example for PUMPKIN. User A cut an inline lemma in an inductive case and proved it using a rewrite:

```
assert (∀ a, S (a + S (a + 0)) = S (S (a + (a + 0)))).
- ... rewrite <- plus_n_0. rewrite -> plus_comm.
```

When we ported User A’s solution to use User B’s definition of `bin_to_nat`, the application of this inline lemma failed. We changed the conclusion of the inline lemma and removed the corresponding rewrite:

```
assert (∀ a, S (a + S a) = S (S (a + a))).
- ... rewrite -> plus_comm.
```

**SEARCH** We used PUMPKIN to search the example for a patch that corresponds to the change in `bin_to_nat`. It found an isomorphism:

```
∀ P b, P (bin_to_nat b) -> P (bin_to_nat b + 0)
∀ P b, P (bin_to_nat b + 0) -> P (bin_to_nat b)
```

**APPLY** After porting to User B’s definition, a rewrite in the proof of the theorem `normalize_correctness` failed:

```
Theorem normalize_correctness:
  ∀ b, nat_to_bin (bin_to_nat b) = normalize b.
Proof.
  ... rewrite -> plus_0_r. X
```

---

<sup>4</sup> If PUMPKIN finds just one implication, it returns that.

Attempting the obvious patch from the difference in tactics—rewriting by `plus_n_0`—failed. Applying the patch that PUMPKIN found fixed the broken proof:

```
... apply patch_inv. rewrite -> plus_0_r.✓
```

In this case, since we ported User A’s definition to a simpler definition,<sup>5</sup> PUMPKIN found a patch that was not the most natural patch. The natural patch would be to remove the `rewrite`, just as we removed a different `rewrite` from the example proof. This did not occur when we ported User B’s definition, which suggests that in the future, a patch finding tool may help inform novice users which definition is simpler: It can factor the proof, then inform the user if two factors are inverses. Tactic-level changes do not provide enough information to determine this; the tool must have a semantic understanding of the terms.

### 3.6.2.1 Configuration

This scenario used the configuration for changes in cases of a fixpoint. Given such a change:

```
Fixpoint f ... := ... | g x
Fixpoint f' ... := ... | g x'
```

PUMPKIN searches two proofs of theorems:

$$\begin{array}{l} \forall \dots, P(f \dots) \\ \forall \dots, P(f' \dots) \end{array}$$

for an isomorphism that corresponds to the change:

$$\begin{array}{l} \forall P, P(x) \rightarrow P(x') \\ \forall P, P(x') \rightarrow P(x) \end{array}$$

The user can apply these patches to fix other broken proofs about the fixpoint.

The key feature that differentiates these from the patches we have encountered so far is that these patches hold for *all*  $P$ ; for changes in fixpoint cases, the procedure abstracts candidates by  $P$ , not by its arguments. PUMPKIN uses this configuration for changes in fixpoint cases:

- 
- 1: *diff* fixpoint cases for goals
  - 2: use *all components* to recursively search an intermediate lemma for a change in conclusions
  - 3: **if** there are candidates **then**
  - 4:   *specialize* and *factor* the candidate  
      *abstract* the factors by functions  
      try to *invert* the patch to find an isomorphism
- 

For the prototype, we require the user to cut the intermediate lemma explicitly and to pass its type and arguments. In the future, an improved semantic differencing component can infer both the

<sup>5</sup> User A uses `*`; User B uses `+`. For arbitrary  $n$ , the term  $2 * n$  reduces to  $n + (n + 0)$ , which does not reduce any further.



<b>Definition</b> divide p q := $\exists$ r, p * r = q.	<b>Definition</b> divide p q := $\exists$ r, q = r * p.
--	--

Figure 6: Old (left) and new (right) definitions of divide in Coq.

intermediate lemma and the arguments: It can search within the proof for some proof of a function that is applied to the fixpoint.

### 3.6.3 Updating Proof Assistant Versions

Coq sometimes makes changes to its standard library that break backwards-compatibility. To test the plausibility of using a patch finding tool for proof assistant version updates, we identified a breaking change in the Coq standard library [51]. The commit changed the definition of `divide` prior to the Coq 8.4 release (Figure 6). The change broke 46 proofs in the standard library. We used PUMPKIN to find an isomorphism that corresponds to the change in `divide`. The isomorphism PUMPKIN found fixed broken proofs.

**REPLAY** We used the proof of `mod_divide` as the example for PUMPKIN. The proof broke with the new `divide`:

```
Theorem mod_divide:
   $\forall$  a b, b $\neq$ 0 -> (a mod b == 0 <-> (divide b a)).
Proof.
... rewrite (div_mod a b Hb) at 2.✗
```

We replayed changes to `mod_divide`:

```
... rewrite mul_comm. symmetry.
rewrite (div_mod a b Hb) at 2.✓
```

**SEARCH** We used PUMPKIN to search the example for a patch that corresponds to the change in `divide`. It found an isomorphism:

```
 $\forall$  r p q, p * r = q -> q = r * p
 $\forall$  r p q, q = r * p -> p * r = q
```

**APPLY** The proof of the theorem `Zmod_divides` broke after rewriting by the changed theorem `mod_divide`:

```
Theorem Zmod_divides:
   $\forall$  a b, b<>0 -> (a mod b = 0 <->  $\exists$  c, a = b * c).
Proof.
... split; intros (c,Hc); exists c; auto.✗
```

Adding the patches PUMPKIN found to a hint database made the proof go through:

```
... split; intros (c,Hc); exists c; auto.✓
```

### 3.6.3.1 Configuration

This scenario used the configuration for changes in dependent arguments to constructors. PUMPKIN searches two proofs that apply the same constructor to different dependent arguments:

$$\begin{array}{c} \dots (C (P \ x)) \dots \\ \dots (C (P' \ x)) \dots \end{array}$$

for an isomorphism between the arguments:

$$\begin{array}{l} \forall x, P \ x \rightarrow P' \ x \\ \forall x, P' \ x \rightarrow P \ x \end{array}$$

The user can apply these patches to patch proofs that apply the constructor (in this case *study*, to fix broken proofs that instantiate *divide* with some specific *r*).

So far, we have encountered changes of this form as arguments to an induction principle; in this case, the change is an argument to a constructor. A patch between arguments to an induction principle maps directly between conclusions of the new and old theorem without induction; a patch between constructors does not. For example, for *divide*, we can find a patch with this form:

$$\forall x, P \ x \rightarrow P' \ x$$

However, without using the induction principle for *exists*, we can't use that patch to prove this:

$$(\exists x, P \ x) \rightarrow (\exists x, P' \ x)$$

This changes the goal type that semantic differencing determines. PUMPKIN uses this configuration for changes in constructor arguments:

- 
- 1: *diff* constructor arguments for goals
  - 2: use *all components* to recursively search those arguments for changes in conclusions
  - 3: if there are candidates **then**
  - 4:   *abstract* the candidate  
      *factor* and try to *invert* the patch to find an isomorphism
- 

For the prototype, the model of constructors for the semantic differencing component is limited, so we ask the user to provide the type of the change in argument (to guide line 2). We can extend semantic differencing to remove this restriction.

## 3.7 CONCLUSION

Rehashing thesis and how we do it

What we haven't accomplished yet at this point (parts of PUMPKIN PATCH future work), segue into next chapter

# 4

---

## PROOF REPAIR ACROSS TYPE EQUIVALENCES

---

This extension to the suite adds support for a broad class of changes in datatypes, handling a large class of practical repair scenarios. What this tool (PUMPKIN Pi) does is, when datatypes change and this breaks a lot of proofs, it generalizes the change in datatype itself (possibly with some user input) so that it can automatically fix proofs broken by the change in datatype.

So in other words, the information from those changes is carried in the difference between the old and new version of the changed datatype, possibly with some user input.

PUMPKIN Pi generalizes that information and applies it automatically.

The work saved is shown on a lot of case studies (see Table from PUMPKIN Pi).

### 4.1 MOTIVATING EXAMPLE

Consider a simple example of using PUMPKIN Pi: repairing proofs after swapping the two constructors of the `list` datatype (Figure 7). This is inspired by a similar change from a user study of proof engineers (Section 4.6). Even such a simple change can cause trouble, as in this proof from the Coq standard library (comments ours for clarity):<sup>1</sup>

```
Lemma rev_app_distr {A} :  
  ∀ (x y : list A), rev (x ++ y) = rev y ++ rev x.
```

---

<sup>1</sup> We use induction instead of pattern matching.

<pre><b>Inductive</b> list (T : <b>Type</b>) :   <b>Type</b> :=     <b>nil</b> : list T     <b>cons</b> : T → list T → list     T.</pre>	<pre><b>Inductive</b> list (T : <b>Type</b>) :   <b>Type</b> :=     <b>cons</b> : T → list T → list     T     <b>nil</b> : list T.</pre>
--	--

Figure 7: A change from the old version of `list` (left) to the new version of `list` (right). The old version of `list` is an inductive datatype that is either empty (the `nil` constructor), or the result of placing an element in front of another `list` (the `cons` constructor). The change swaps these two constructors (orange).

```

swap T (l : Old.list T) : New. swap-1 T (l : New.list T) :
  list T := Old.list T :=
  Old.list_rect T (fun (l : Old.list T) => New.list T)
    )
  New.nil
  (fun t _ (IH1 : New.list T) => New.cons T t IH1)
  1.
  New.list_rect T (fun (l : New.list T) => Old.list T)
    )
  (fun t _ (IH1 : Old.list T) => Old.cons T t IH1)
  Old.nil
  1.

Lemma section: ∀ T (l : Old. list T),
  swap-1 T (swap T l) = 1.
Proof.
  intros T l. symmetry.
  induction l as [ | a l0 H ].
  - auto.
  - simpl. rewrite ← H. auto.
Qed.

Lemma retraction: ∀ T (l : New. list T),
  swap T (swap-1 T l) = 1.
Proof.
  intros T l. symmetry.
  induction l as [t | t0 H].
  - simpl. rewrite ← H. auto.
  - auto.
Qed.

```

Figure 8: Two functions between `Old.list` and `New.list` (top) that form an equivalence (bottom).

```

Proof. (* by induction over x and y *)
  induction x as [| a l IH1].
  (* x nil: *) induction y as [| a l IH1].
  (* y nil: *) simpl. auto.
  (* y cons *) simpl. rewrite app_nil_r; auto.
  (* both cons: *) intro y. simpl.
  rewrite (IH1 y). rewrite app_assoc; trivial.
Qed.

```

This lemma says that appending (`++`) two lists and reversing (`rev`) the result behaves the same as appending the reverse of the second list onto the reverse of the first list. The proof script works by induction over the input lists `x` and `y`: In the base case for both `x` and `y`, the result holds by reflexivity. In the base case for `x` and the inductive case for `y`, the result follows from the existing lemma `app_nil_r`. Finally, in the inductive case for both `x` and `y`, the result follows by the inductive hypothesis and the existing lemma `app_assoc`.

When we change the list type, this proof no longer works. To repair this proof with PUMPKIN Pi, we run this command:

```
Repair Old.list New.list in rev_app_distr.
```

assuming the old and new list types from Figure 7 are in modules `Old` and `New`. This suggests a proof script that succeeds (in light blue to denote PUMPKIN Pi produces it automatically):

```

Proof. (* by induction over x and y *)
  intros x. induction x as [a l IH1]; intro y0.
  - (* both cons: *) simpl. rewrite IH1. simpl.
    rewrite app_assoc. auto.
  - (* x nil: *) induction y0 as [a l H].
    + (* y cons: *) simpl. rewrite app_nil_r. auto.
    + (* y nil: *) auto.
Qed.

```

where the dependencies (`rev`, `++`, `app_assoc`, and `app_nil_r`) have also been updated automatically ①. If we would like, we can manually modify this to something that more closely matches the style of the original proof script:

```
Proof. (* by induction over x and y *)
  induction x as [a l IHl|].
  (* both cons: *) intro y. simpl.
  rewrite (IHl y). rewrite app_assoc; trivial.
  (* x nil: *) induction y as [a l IHl|].
  (* y cons: *) simpl. rewrite app_nil_r; auto.
  (* y nil: *) simpl. auto.
Qed.
```

We can even repair the entire list module from the Coq standard library all at once by running the `Repair module` command ①. When we are done, we can get rid of `Old.list`.

The key to success is taking advantage of Coq’s structured proof term language: Coq compiles every proof script to a proof term in a rich functional programming language called Gallina—PUMPKIN Pi repairs that term. PUMPKIN Pi then decompiles the repaired proof term (with optional hints from the original proof script) back to a suggested proof script that the proof engineer can maintain.

In contrast, updating the poorly structured proof script directly would not be straightforward. Even for the simple proof script above, grouping tactics by line, there are  $6! = 720$  permutations of this proof script. It is not clear which lines to swap since these tactics do not have a semantics beyond the searches their evaluation performs. Furthermore, just swapping lines is not enough: even for such a simple change, we must also swap arguments, so `induction x as [| a l IHl ]` becomes `induction x as [a l IHl|]`. A recent thesis [74] describes the challenges of repairing tactics in detail. PUMPKIN Pi’s approach circumvents this challenge.

## 4.2 APPROACH

PUMPKIN Pi can do much more than permute constructors. Given an equivalence between types  $A$  and  $B$ , PUMPKIN Pi repairs functions and proofs defined over  $A$  to instead refer to  $B$ . It does this in a way that allows for removing references to  $A$ , which is essential for proof repair, since  $A$  may be an old version of an updated type.

Like I mentioned earlier, this also works using differencing and program transformations over proof terms. Here, differencing thus looks at the difference between versions of the changed datatype, and finds something called a type equivalence (Section 4.2.1). Sometimes differencing is automatic, and sometimes it’s manual. Then, program transformation ports proofs across the equivalence directly (Section 4.2.2). So they take care of application.

Figure 9 shows how this comes together when the proof engineer invokes PUMPKIN Pi:

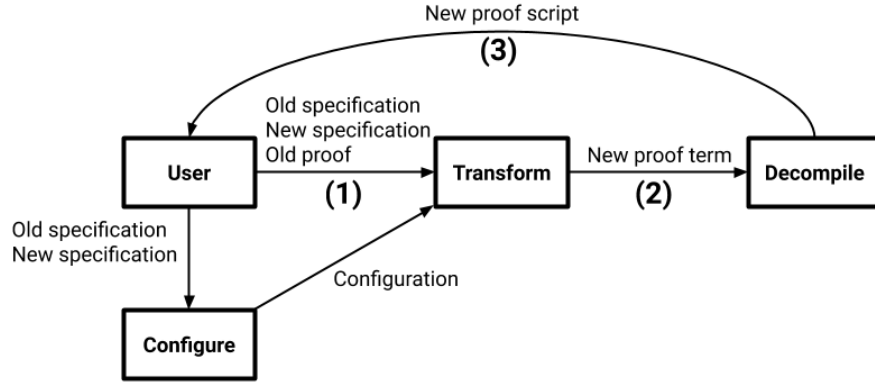


Figure 9: The workflow for PUMPKIN Pi.

1. The proof engineer **Configures** PUMPKIN Pi, either manually or automatically.
2. The configured **Transform** transforms the old proof term into the new proof term.
3. **Decompile** suggests a new proof script.

There are currently four search procedures for automatic configuration implemented in PUMPKIN Pi (see Table 1 on page 43). Manual configuration makes it possible for the proof engineer to configure the transformation to any equivalence, even without a search procedure. Section 4.6 shows examples of both workflows applied to real scenarios

#### 4.2.1 Scope: Type Equivalences

PUMPKIN Pi repairs proofs in response to changes in types that correspond to *type equivalences* [80], or pairs of functions that map between two types and are mutual inverses.<sup>2</sup> When a type equivalence between types  $A$  and  $B$  exists, those types are *equivalent* (denoted  $A \simeq B$ ). Figure 8 shows a type equivalence between the two versions of `list` from Figure 7 that PUMPKIN Pi discovered and proved automatically ①.

To give some intuition for what kinds of changes can be described by equivalences, we preview two changes below. See Table 1 on page 43 for more examples.

**Factoring out Constructors.** Consider changing the type  $I$  to the type  $J$  in Figure 10.  $J$  can be viewed as  $I$  with its two constructors  $A$  and  $B$  pulled out to a new argument of type `bool` for a single constructor. With PUMPKIN Pi, the proof engineer can repair functions and proofs about  $I$  to instead use  $J$ , as long as she configures PUMPKIN Pi to describe which constructor of  $I$  maps to `true` and which maps to `false`. This information about constructor mappings induces an equivalence

<sup>2</sup> The adjoint follows, and PUMPKIN Pi includes machinery to prove it ⑩ ②③.

$I \simeq J$  across which PUMPKIN Pi repairs functions and proofs. File ② shows an example of this, mapping  $A$  to true and  $B$  to false, and repairing proofs of De Morgan’s laws.

**Adding a Dependent Index.** At first glance, the word *equivalence* may seem to imply that PUMPKIN Pi can support only changes in which the proof engineer does not add or remove information. But equivalences are more powerful than they may seem. Consider, for example, changing a list to a length-indexed vector (Figure 11). PUMPKIN Pi can repair functions and proofs about lists to functions and proofs about vectors of particular lengths ③, since  $\Sigma(l : \text{list } T). \text{length } l = n \simeq \text{vector } T \ n$ . From the proof engineer’s perspective, after updating specifications from `list` to `vector`, to fix her functions and proofs, she must additionally prove invariants about the lengths of her lists. PUMPKIN Pi makes it easy to separate out that proof obligation, then automates the rest.

More generally, in homotopy type theory, with the help of quotient types, it is possible to form an equivalence from a relation, even when the relation is not an equivalence [3]. While Coq lacks quotient types, it is possible to achieve a similar outcome and use PUMPKIN Pi for changes that add or remove information when those changes can be expressed as equivalences between  $\Sigma$  types or sum types.

#### 4.2.2 Goal: Transport with a Twist

The goal of PUMPKIN Pi is to implement a kind of proof reuse known as *transport* [80], but in a way that is suitable for repair. Informally, transport takes a term  $t$  and produces a term  $t'$  that is the same as  $t$  modulo an equivalence  $A \simeq B$ . If  $t$  is a function, then  $t'$  behaves the same way modulo the equivalence; if  $t$  is a proof, then  $t'$  proves the same theorem the same way modulo the equivalence.

When transport across  $A \simeq B$  takes  $t$  to  $t'$ , we say that  $t$  and  $t'$  are *equal up to transport* across that equivalence (denoted  $t \equiv_{A \simeq B} t'$ ).<sup>3</sup> In Section 4.1, the original `append` function `++` over `Old.list` and the repaired `append` function `++` over `New.list` that PUMPKIN Pi produces are equal up to transport across the equivalence from Figure 8, since (by `app_ok` ①):

$$\forall T (l1 \ l2 : \text{Old.list } T), \\ \text{swap } T (l1 ++ l2) = (\text{swap } T \ l1) ++ (\text{swap } T \ l2).$$

The original `rev_app_distr` is equal to the repaired proof up to transport, since both prove the same thing the same way up to the equivalence, and up to the changes in `++` and `rev`.

<sup>3</sup> This notation should be interpreted in a metatheory with *univalence*—a property that Coq lacks—or it should be approximated in Coq. The details of transport with univalence are in the Homotopy Type Theory book [80], and an approximation in Coq is in the univalent parametricity framework paper [77]. For equivalent  $A$  and  $B$ , there can be many equivalences  $A \simeq B$ . Equality up to transport is across a *particular* equivalence, but we erase this in the notation.

```

Inductive I :=
| A : I
| B : I.

Inductive J :=
| makeJ : bool → J.

```

Figure 10: The old type  $I$  (left) is either  $A$  or  $B$ . The new type  $J$  (right) is  $I$  with  $A$  and  $B$  factored out to `bool` (orange).

```

Inductive list (T : Type) :
  Type :=
| nil : list T
| cons : T → list T → list T.

Inductive vector (T : Type) : nat →
  Type :=
| nil : vector T 0
| cons : T → ∀ (n : nat), vector T n → vector T (S n).

```

Figure 11: A vector (bottom) is a list (top) indexed by its length (orange). Vectors effectively make it possible to enforce length invariants about lists at compile time.

Transport typically works by applying the functions that make up the equivalence to convert inputs and outputs between types. This approach would not be suitable for repair, since it does not make it possible to remove the old type  $A$ . PUMPKIN Pi implements transport in a way that allows for removing references to  $A$ —by proof term transformation.

### 4.3 DIFFERENCING

General idea from PUMPKIN Pi, explain automatic configurations and so on, segue into DEVoid example

#### 4.3.1 Algebraic Ornaments

An algebraic ornament relates an inductive type  $A$  to an indexed version of that type  $B$  with a new index of type  $I_B$ , where the new index is fully determined by a unique fold over  $A$ . For example, `vector` is exactly `list` with a new index of type `nat`, where the new index is fully determined by the `length` function. Consequentially, there are two functions:

```

ltv : list T → Σ(n : nat).vector T n.
vtl : Σ(n : nat).vector T n → list T.

```

that are mutual inverses:

```

∀ (l : list T),          vtl (ltv l) = l.
∀ (v : Σ(n : nat).vector T n), ltv (vtl v) = v.

```

and therefore form the type equivalence from Section ?? . Moreover, since the new index is fully determined by `length`, we can relate `length` to `ltv`:

```

∀ (l : list T), length l = πl (ltv l).

```

In general, we can view an algebraic ornament as a type equivalence:

$$A \vec{i} \simeq \Sigma(n : I_B \vec{i}). B \text{ (index } n \vec{i})$$



where  $\vec{i}$  are the indices of  $A$ ,  $I_B$  is a function over those indices, and the index operation inserts the new index  $n$  at the right offset. Such a type equivalence consists of two functions [80]:

$$\begin{aligned} \text{promote} & : A \ \vec{i} \rightarrow \Sigma(n : I_B \ \vec{i}).B \ (\text{index } n \ \vec{i}). \\ \text{forget} & : \Sigma(n : I_B \ \vec{i}).B \ (\text{index } n \ \vec{i}) \rightarrow A \ \vec{i}. \end{aligned}$$

that are mutual inverses:<sup>4</sup>

$$\begin{aligned} \text{section} & : \forall (a : A \ \vec{i}), \quad \text{forget} \ (\text{promote } a) = a. \\ \text{retraction} & : \forall (b_\Sigma : \Sigma(n : I_B \ \vec{i}).B \ (\text{index } n \ \vec{i})), \text{promote} \ (\text{forget } b_\Sigma) = b_\Sigma. \end{aligned}$$

An algebraic ornament is additionally equipped with an indexer, which is a unique fold:

$$\text{indexer} : A \ \vec{i} \rightarrow I_B \ \vec{i}.$$

which projects the promoted index:

$$\text{coherence} : \forall (a : A \ \vec{i}), \text{indexer } a = \pi_l \ (\text{promote } a).$$

Following existing work [46], we call this equivalence the *ornamental promotion isomorphism*; when it holds and the indexer exists, we say that  $B$  is an algebraic ornament of  $A$ .

`Find ornament` searches for algebraic ornaments between types and is, to the best of our knowledge, the first search algorithm for ornaments.

In their original form, ornaments are a programming mechanism: Given a type  $A$ , an ornament determines some new type  $B$ . We invert this process for algebraic ornaments: Given types  $A$  and  $B$ , `DEVoid` searches for an ornament between them. This is possible for algebraic ornaments precisely because the indexer is extensionally unique. For example, all possible indexers for `list` and `vector` must compute the length of a list; if we were to try doubling the length instead, we would not be able to satisfy the equivalence.

`Find ornament` takes two inductive types and searches for the components of the ornamental promotion isomorphism between them:

- **Inputs:** Inductive types  $A$  and  $B$ , assuming:
  - $B$  is an algebraic ornament of  $A$ ,
  - $B$  has the same number of constructors in the same order as  $A$ ,
  - $A$  and  $B$  do not contain recursive references to themselves under products, and
  - for every recursive reference to  $A$  in  $A$ , there is exactly one new hypothesis in  $B$ , which is exactly the new index of the corresponding recursive reference in  $B$ .
- **Outputs:** Functions `promote`, `forget`, and `indexer`, guaranteeing:

<sup>4</sup> The adjunction condition follows from section and retraction.

- the outputs form the ornamental promotion isomorphism between the inputs.

`Find ornament` includes an option to generate a proof that the outputs form the ornamental promotion isomorphism; by default, this option is false, since `Lift` does not need this proof.

**Presentation.** We present both algorithms relationally, using a set of judgments; to turn these relations into algorithms, prioritize the rules by running the derivations in order, falling back to the original term when no rules match. The default rule for a list of terms is to run the derivation on each element of the list individually.

**Notes on Syntax.** The language the algorithms operate over is  $\text{CIC}_\omega$  with primitive eliminators; this is a simplified version of the type theory underlying Coq. Figure 12 contains the syntax (which includes variables, sorts, product types, functions, inductive types, constructors, and eliminators), as well as the syntax for some judgments and operations, the rules for which are standard and thus omitted. For simplicity of presentation, we assume variables are names; we assume that all names are fresh. As in Coq, we assume the existence of an inductive type  $\Sigma$  for sigma types with projections  $\pi_l$  and  $\pi_r$ ; for simplicity, we assume projections are primitive. Throughout, we use  $\vec{i}$  and  $\{t_1, \dots, t_n\}$  to denote lists of terms, and we use  $\vec{i}[j]$  to denote accessing the element of the list  $\vec{i}$  at offset  $j$ .

**Common Definitions.** The algorithms assume list insertion and removal functions `insert` and `remove`, plus two functions `DEVOID` implements: `off` computes the offset of the new index of type  $I_B$  in  $B$ 's indices, and `new` determines whether a hypothesis in a case of the eliminator type of  $B$  is new. Figure 13 contains other common definitions, the names for which are reserved: The `index` and `deindex` functions insert an index into and remove an index from a list at the index computed by `off`. Input type  $A$  expands to an inductive type with indices of types  $\vec{X}_A$ , sort  $s_A$ , and constructors  $\{C_{A_1}, \dots, C_{A_n}\}$ .  $P_A$  denotes the type of the motive of the eliminator of  $A$ , and each  $E_{A_i}$  denotes the type of the eliminator for the  $i$ th constructor of  $A$ . Analogous names are also reserved for input type  $B$ .

The `Find ornament` algorithm implements the specification. It builds on three intermediate steps: one to generate each of `indexer`, `promote`, and `forget`. Figure 14 shows the algorithm for generating `indexer`. The algorithms for generating `promote` and `forget` are similar; Figure 15 shows only the derivations for generating `promote` that are different from those for generating `indexer`, and the derivations for generating `forget` are omitted.

$\langle i \rangle \in \mathbb{N}, \langle v \rangle \in \text{Vars}, \langle s \rangle \in \{ \text{Prop}, \text{Set} \}, \Gamma \vdash t : T$  // type checking  
 $\text{Type}(\langle i \rangle)$   $\Gamma \vdash t_1 \equiv_{\beta\delta_i} t_2$  // definitional equality  
 $\langle t \rangle ::= \langle v \rangle \mid \langle s \rangle \mid \Pi(\langle v \rangle : \langle t \rangle). \langle t \rangle \mid t_\beta$  // beta-reduction  
 $\lambda(\langle v \rangle : \langle t \rangle). \langle t \rangle \mid \langle t \rangle \langle t \rangle \mid t_{\beta\delta_i}$  // normalization  
 $\text{Ind}(\langle v \rangle : \langle t \rangle)\{\langle t \rangle, \dots, \langle t \rangle\} \mid \text{Constr } t [y / x]$  // substitution  
 $(\langle i \rangle, \langle t \rangle) \mid \xi(I, Q, c, C)$  // type of  
 $\text{Elim}(\langle t \rangle, \langle t \rangle)\{\langle t \rangle, \dots, \langle t \rangle\}$  // eliminator

Figure 12:  $\text{CIC}_\omega$  syntax (left, from existing work [79]) and judgments and operations (right).

$A := \text{Ind}(\text{Ty}_A : \Pi(\vec{i}_A : \vec{X}_A).s_A)\{C_{A_1}, \dots, C_{A_n}\}$   
 $B := \text{Ind}(\text{Ty}_B : \Pi(\vec{i}_B : \vec{X}_B).s_B)\{C_{B_1}, \dots, C_{B_n}\}$   
 $\forall 1 \leq i \leq n,$   
 $E_{A_i}(p_A : P_A) := \xi(A, p_A, \text{Constr}(i, A), C_{A_i})$   
 $E_{B_i}(p_B : P_B) := \xi(B, p_B, \text{Constr}(i, B), C_{B_i})$   
 $P_A := \Pi(\vec{i}_A : \vec{X}_A)(a : A \vec{i}_A).s_A$   
 $P_B := \Pi(\vec{i}_B : \vec{X}_B)(b : B \vec{i}_B).s_B$   
 $\text{index} := \text{insert}(\text{off } A \ B)$   
 $\text{deindex} := \text{remove}(\text{off } A \ B)$

Figure 13: Common definitions for both algorithms.

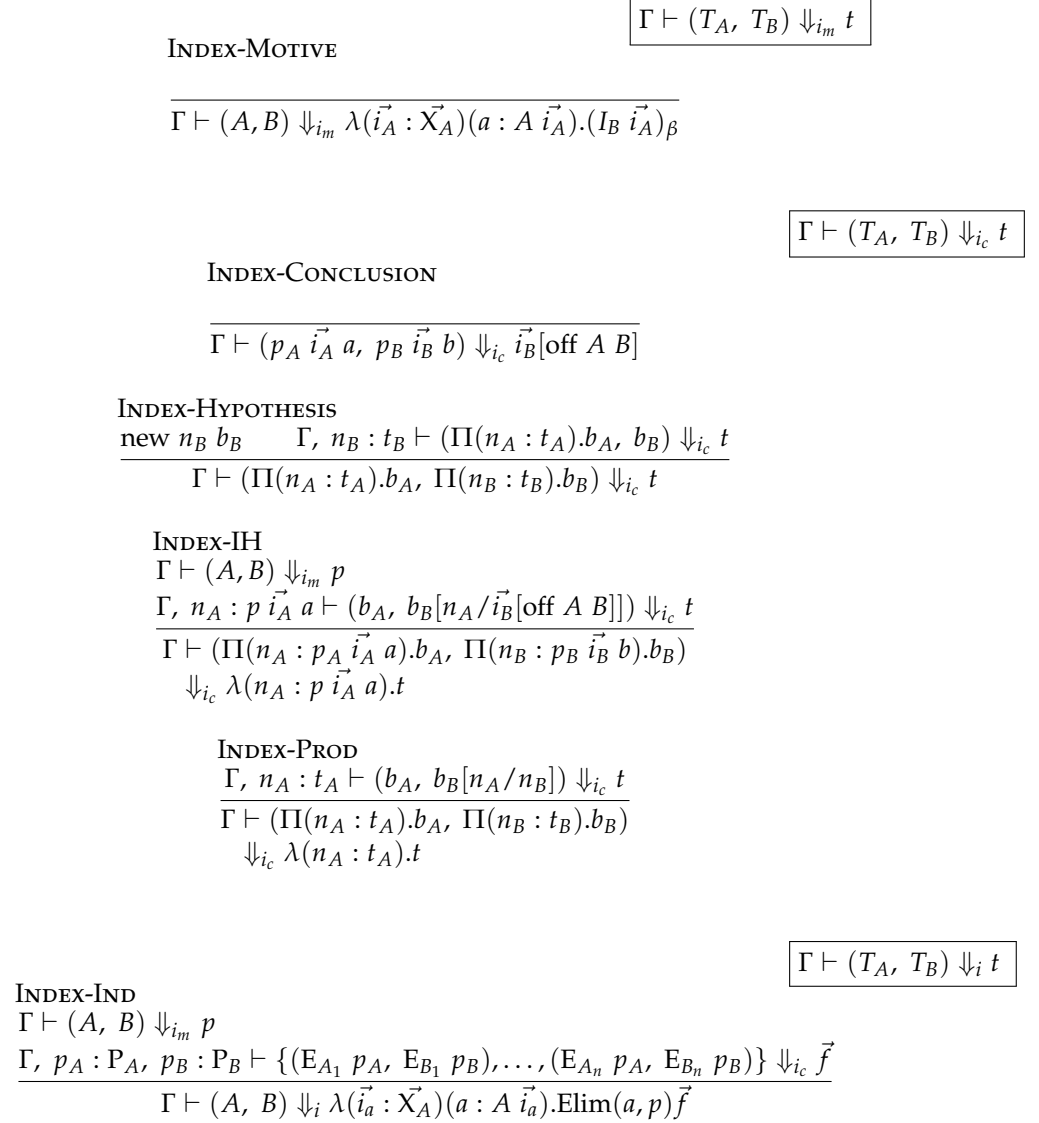


Figure 14: Identifying the indexer function.

#### 4.3.1.1 Searching for the Indexer

Search generates the `indexer` by traversing the types of the eliminators for  $A$  and  $B$  in parallel using the algorithm from Figure 14, which consists of three judgments: one to generate the motive, one to generate each case, and one to compose the motive and cases.

**Generating the Motive.** The  $(T_A, T_B) \Downarrow_{i_m} t$  judgment consists of only the derivation **INDEX-MOTIVE**, which computes the indexer motive from the types  $A$  and  $B$  (expanded in Figure 13). It does this by constructing a function with  $A$  and its indices as premises, and the type  $I_B$  in the conclusion with the appropriate indices. Consider `list` and `vector`:

```
list T := Ind (TyA : Type) {...}      vector T := Ind (TyB : Π
  (n : nat).Type) {...}
```

For these types, INDEX-MOTIVE computes the motive:

```
λ (l:list T) . nat
```

**Generating Each Case.** The  $\Gamma \vdash (T_A, T_B) \Downarrow_{i_c} t$  judgment generates each case of the indexer by traversing in parallel the corresponding cases of the eliminator types for  $A$  and  $B$ . It consists of four derivations: INDEX-CONCLUSION handles base cases and conclusions of inductive cases, while INDEX-HYPOTHESIS, INDEX-IH, and INDEX-PROD recurse into products.

INDEX-HYPOTHESIS handles each new hypothesis that corresponds to a new index in an inductive hypothesis of an inductive case of the eliminator type for  $B$ . It adds the new index to the environment, then recurses into the body of only the type for which the index already exists. For example, in the inductive case of `list` and `vector`, new determines that `n` is the new hypothesis. INDEX-HYPOTHESIS then recurses into the body of only the `vector` case:

```
Π (tl:T) (l:list T) (IHl:pA 1), ...      Π (tv:T) (v:vector T n)
  (IHv:pB n v), ...
```

INDEX-PROD is next. It recurses into product types when the hypothesis is neither a new index nor an inductive hypothesis. Here, it runs twice, recursing into the body and substituting names until it hits the inductive hypothesis for both types:

```
Π (IHl:pA 1), pA (cons tl 1)      Π (IHv:pB n 1), pB
  (S n) (consV n tl 1)
```

INDEX-IH then takes over. It substitutes the new motive in the inductive hypothesis, then recurses into both bodies, substituting the new inductive hypothesis for the index in the eliminator type for  $B$ . Here, it substitutes the new motive for  $p_A$  in the type of  $IH_l$ , extends the environment with  $IH_l$ , then substitutes  $IH_l$  for  $n$ , so that it recurses on these types:

```
pA (cons tl 1)      pB (S IHl) (consV IHl tl 1)
```

Finally, INDEX-CONCLUSION computes the conclusion by taking the index of motive  $p_B$  at off  $A B$ , here  $S IH_l$ . In total, this produces a function that computes the length of `cons t 1`:

```
λ (tl:T) (l:list T) (IHl:(λ (l:list T).nat) 1).S IHl
```

**Composing the Result.** The  $\Gamma \vdash (T_A, T_B) \Downarrow_i t$  judgment consists of only INDEX-IND, which identifies the motive and each case using the other two judgments, then composes the result. In the case of `list` and `vector`, this produces a function that computes the length of a list:

```
λ (l:list T).Elim(1, λ (l:list T).nat)
  {0, λ (tl:T) (l:list T) (IHl:(λ (l:list T).nat) 1).S IHl}
```

$$\begin{array}{c}
\boxed{\Gamma \vdash (T_A, T_B) \Downarrow_{p_m} t} \\
\text{PROMOTE-MOTIVE} \\
\frac{\Gamma \vdash (A, B) \Downarrow_i \pi}{\Gamma \vdash (A, B) \Downarrow_{p_m} \lambda(\vec{i}_a : \vec{X}_A)(a : A \vec{i}_a).B \text{ (index } (\pi \vec{i}_a a) \vec{i}_a)} \\
\\
\boxed{\Gamma \vdash (T_A, T_B) \Downarrow_{p_c} t} \\
\text{PROMOTE-CONCLUSION} \\
\frac{}{\Gamma \vdash (p_A \vec{i}_A a, p_B \vec{i}_B b) \Downarrow_{p_c} b} \\
\\
\text{PROMOTE-IH} \\
\frac{\Gamma \vdash (A, B) \Downarrow_i \pi \quad \Gamma \vdash (A, B) \Downarrow_{p_m} p \quad \Gamma, n_A : p \vec{i}_A a \vdash (b_A, b_B[n_A/b][\pi \vec{i}_A a / \vec{i}_B[\text{off } A \ B]]) \Downarrow_{p_c} t}{\Gamma \vdash (\Pi(n_A : p_A \vec{i}_A a).b_A, \Pi(n_B : p_B \vec{i}_B b).b_B) \Downarrow_{p_c} \lambda(n_A : p \vec{i}_A a).t} \\
\\
\boxed{\Gamma \vdash (T_A, T_B) \Downarrow_p t} \\
\text{PROMOTE-IND} \\
\frac{\Gamma \vdash (A, B) \Downarrow_i \pi \quad \Gamma \vdash (A, B) \Downarrow_{p_m} p \quad \Gamma, p_A : P_A, p_B : P_B \vdash \{(E_{A_1} p_A, E_{B_1} p_B), \dots, (E_{A_n} p_A, E_{B_n} p_B)\} \Downarrow_{p_c} \vec{f}}{\Gamma \vdash (A, B) \Downarrow_p \lambda(\vec{i}_A : \vec{X}_A)(a : A \vec{i}_A).\exists (\pi \vec{i}_A a) (\text{Elim}(a, p)\vec{f})}
\end{array}$$

Figure 15: Identifying the promotion function.

#### 4.3.1.2 Searching for Promote and Forget

Figure 15 shows the interesting derivations for the judgment  $(T_A, T_B) \Downarrow_p t$  that searches for promote: PROMOTE-MOTIVE identifies the motive as  $B$  with a new index (which it computes using `indexer`, denoted by metavariable  $\pi$ ). When PROMOTE-IH recurses, it substitutes the inductive hypothesis for the term rather than for its index, and it substitutes the new index (which it also computes using `indexer`) inside of that term. PROMOTE-CONCLUSION returns the entire term, rather than its index. Finally, PROMOTE-IND not only recurses into each case, but also packs the result.

The omitted derivations to search for forget are similar, except that the domain and range are switched. Consequentially, `indexer` is never needed; FORGET-MOTIVE removes the index rather than inserting it, and FORGET-IH no longer substitutes the index. Additionally, FORGET-HYPOTHESIS adds the hypothesis for the new index rather than skipping it, and FORGET-IND eliminates over the projection rather than packing the result.

#### 4.3.1.3 Core Search Algorithm

The core search algorithm produces `indexer`, `promote`, and `forget`, then composes them into a tuple. This tuple is how DEVOLD represents

ornaments internally. DEVOID has options (used in `Example.v`) that tell search to generate proofs that its outputs are correct, thereby increasing confidence in and usefulness of those outputs. The proof of coherence is reflexivity. The intuition behind the automation to prove section and retraction (`equivalence.ml`) is that `promote` and `forget` map along corresponding constructors, so inductive cases preserve equalities. Thus, each inductive case of these proofs is generated by a fold that rewrites each recursive reference, with reflexivity as identity.

#### 4.3.2 *Other Search Procedures*

Brief explanation of these and how differencing works for them in detail based on algebraic ornaments

#### 4.3.3 *Designing New Search Procedures*

How hard, how useful

Limitations and whether they're addressed in other tools yet or not

### 4.4 TRANSFORMATION

Parts of PUMPKIN Pi Transformation, with DEVOID 3.2 and 4.2 as examples, plus some of the beautiful Carlo theory to explain why we go from equivalences to configurations and what that really means

How the transformation works in detail

Limitations and whether they're addressed in other tools yet or not

### 4.5 IMPLEMENTATION

Parts of PUMPKIN Pi and DEVOID implementation, plus more (still need to arrange, fill in, and so on)

#### 4.5.1 *Tool Details*

Implemented in blah blah blah, and so on.

#### 4.5.2 *Workflow Integration*

##### 4.5.2.1 *Configure*

##### 4.5.2.2 *Transform*

**Termination.** When a subterm unifies with a configuration term, this suggests that PUMPKIN Pi *can* transform the subterm, but it does not necessarily mean that it *should*. In some cases, doing so would result in nontermination. For example, if  $B$  is a refinement of  $A$ , then we

$$\langle v \rangle \in \text{Vars}, \langle t \rangle \in \text{CIC}_\omega$$

$$\langle p \rangle ::= \text{intro } \langle v \rangle \mid \text{rewrite } \langle t \rangle \langle t \rangle \mid \text{symmetry} \mid \text{apply } \langle t \rangle \mid \text{induction} \\ \langle t \rangle \langle t \rangle \{ \langle p \rangle, \dots, \langle p \rangle \} \mid \text{split } \{ \langle p \rangle, \langle p \rangle \} \mid \text{left} \mid \text{right} \mid \langle p \rangle . \langle p \rangle$$

Figure 16: Qtac syntax.

can always run EQUIVALENCE over and over again, forever. We thus include some simple termination checks in our code (12).

**Intent.** Even when termination is guaranteed, whether to transform a subterm depends on intent. That is, PUMPKIN Pi automates the case of porting *every*  $A$  to  $B$ , but proof engineers sometimes wish to port only *some*  $A$ s to  $B$ s. PUMPKIN Pi has some support for this using an interactive workflow (13), with plans for automatic support in the future.

**From  $\text{CIC}_\omega$  to Coq.** The implementation (4) of the transformation handles language differences to scale from  $\text{CIC}_\omega$  to Coq. We use the existing Preprocess [73] command to turn pattern matching and fixpoints into eliminators. We handle refolding of constants in constructors using DepConstr.

**Reaching Real Proof Engineers.** Many of our design decisions in implementing PUMPKIN Pi were informed by our partnership with an industrial proof engineer (Section 4.6). For example, the proof engineer rarely had the patience to wait more than ten seconds for PUMPKIN Pi to port a term, so we implemented optional aggressive caching, even caching intermediate subterms encountered while running the transformation (14). We also added a cache to tell PUMPKIN Pi not to  $\delta$ -reduce certain terms (14). With these caches, the proof engineer found PUMPKIN Pi efficient enough to use on a code base with tens of thousands of lines of code and proof.

The experiences of proof engineers also inspired new features. For example, we implemented a search procedure to generate custom eliminators to help reason about types like  $\Sigma(l : \text{list } T).\text{length } l = n$  by reasoning separately about the projections (15). We added informative error messages (22) to help the proof engineer distinguish between user errors and bugs. These features helped with workflow integration.

#### 4.5.2.3 Decompile

**Transform** produces a proof term, while the proof engineer typically writes and maintains proof scripts made up of tactics. We improve usability thanks to the realization that, since Coq’s proof term language Gallina is very structured, we can decompile these Gallina terms to suggested Ltac proof scripts for the proof engineer to maintain.

**Decompile** implements a prototype of this translation (11): it translates a proof term to a suggested proof script that attempts to prove



$$\begin{array}{c}
\boxed{\Gamma \vdash t \Rightarrow p} \\
\\
\text{INTRO} \quad \frac{\Gamma, n : T \vdash b \Rightarrow p}{\Gamma \vdash \lambda(n : T).b \Rightarrow \text{intro } n. p} \qquad \text{SYMMETRY} \quad \frac{\Gamma \vdash H \Rightarrow p}{\Gamma \vdash \text{eq\_sym } H \Rightarrow \text{symmetry. } p} \\
\\
\text{SPLIT} \quad \frac{\Gamma \vdash l \Rightarrow p \quad \Gamma \vdash r \Rightarrow q}{\Gamma \vdash \text{Constr}(0, \wedge) l r \Rightarrow \text{split}\{p, q\}.} \\
\\
\text{LEFT} \quad \frac{\Gamma \vdash H \Rightarrow p}{\Gamma \vdash \text{Constr}(0, \vee) H \Rightarrow \text{left. } p} \qquad \text{RIGHT} \quad \frac{\Gamma \vdash H \Rightarrow p}{\Gamma \vdash \text{Constr}(1, \vee) H \Rightarrow \text{right. } p} \\
\\
\text{REWRITE} \quad \frac{\Gamma \vdash H_1 : x = y \quad \Gamma \vdash H_2 \Rightarrow p}{\Gamma \vdash \text{Elim}(H_1, P)\{x, H_2, y\} \Rightarrow \text{symmetry. rewrite } P H_1. p} \\
\\
\text{INDUCTION} \quad \frac{\Gamma \vdash \vec{f} \Rightarrow \vec{p}}{\Gamma \vdash \text{Elim}(t, P) \vec{f} \Rightarrow \text{induction } P t \vec{p}} \qquad \text{APPLY} \quad \frac{\Gamma \vdash t \Rightarrow p}{\Gamma \vdash ft \Rightarrow \text{apply } f. p} \\
\\
\text{BASE} \quad \frac{}{\Gamma \vdash t \Rightarrow \text{apply } t}
\end{array}$$

Figure 17: Qtac decompiler semantics.

the same theorem the same way. Note that this problem is not well defined: while there is always a proof script that works (applying the proof term with the `apply` tactic), the result is often qualitatively unreadable. This is the baseline behavior to which the decompiler defaults. The goal of the decompiler is to improve on that baseline as much as possible, or else suggest a proof script that is close enough to correct that the proof engineer can manually massage it into something that works and is maintainable.

**Decompile** achieves this in two passes: The first pass decompiles proof terms to proof scripts that use a predefined set of tactics. The second pass improves on suggested tactics by simplifying arguments, substituting tacticals, and using hints like custom tactics and decision procedures.

**First Pass: Basic Proof Scripts.** The first pass takes Coq terms and produces tactics in Ltac, the proof script language for Coq. Ltac can be confusing to reason about, since Ltac tactics can refer to Gallina terms, and the semantics of Ltac depends both on the semantics of Gallina and on the implementation of proof search procedures written in OCaml. To give a sense of how the first pass works without the clutter of these details, we start by defining a mini decompiler that implements a simplified version of the first pass. Section ?? explains how we scale this to the implementation.

The mini decompiler takes  $CIC_\omega$  terms and produces tactics in a mini version of Ltac which we call Qtac. The syntax for Qtac is in Figure 16. Qtac includes hypothesis introduction (`intro`), rewriting (`rewrite`), symmetry of equality (`symmetry`), application of a term to prove the goal (`apply`), induction (`induction`), case splitting of conjunctions (`split`), constructors of disjunctions (`left` and `right`), and composition (`.`). Unlike in Ltac, `induction` and `rewrite` take a motive explicitly (rather than relying on unification), and `apply` creates a new subgoal for each function argument.

The semantics for the mini decompiler  $\Gamma \vdash t \Rightarrow p$  are in Figure 17 (assuming  $=$ , `eq_sym`,  $\wedge$ , and  $\vee$  are defined as in Coq). As with the real decompiler, the mini decompiler defaults to the proof script that applies the entire proof term with `apply` (`BASE`). Otherwise, it improves on that behavior by recursing over the proof term and constructing a proof script using a predefined set of tactics.

For the mini decompiler, this is straightforward: Lambda terms become introduction (`INTRO`). Applications of `eq_sym` become symmetry of equality (`SYMMETRY`). Constructors of conjunction and disjunction map to the respective tactics (`SPLIT`, `LEFT`, and `RIGHT`). Applications of equality eliminators compose symmetry (to orient the rewrite direction) with rewrites (`REWRITE`), and all other applications of eliminators become induction (`INDUCTION`). The remaining applications become apply tactics (`APPLY`). In all cases, the decompiler recurses, breaking into cases, until only the `BASE` case holds.

```

fun (y0 : list A) =>
  list_rect _ _ (fun a l H =>
    eq_ind_r _ eq_refl (app_nil_r (rev l) (a::[])))
    eq_refl
    y0

- intro y0. induction y0 as [a l H].
+ simpl. rewrite app_nil_r. auto.
+ auto.

```

Figure 18: Proof term (top) and decompiled proof script (bottom) for the base case of `rev_app_distr` (Section 4.1), with corresponding terms and tactics grouped by color & number.

While the mini decompiler is very simple, only a few small changes are needed to move this to Coq. The generated proof term of `rev_app_distr` from Section 4.1, for example, consists only of induction, rewriting, simplification, and reflexivity (solved by `auto`). Figure 18 shows the proof term for the base case of `rev_app_distr` alongside the proof script that PUMPKIN Pi suggests. This script is fairly low-level and close to the proof term, but it is already something that the proof engineer can step through to understand, modify, and maintain. There are few differences from the mini decompiler needed to produce this, for example handling of rewrites in both directions (`eq_ind_r` as opposed to `eq_ind`), simplifying rewrites, and turning applications of `eq_refl` into reflexivity or `auto`.

**Second Pass: Better Proof Scripts.** The implementation of **Decompile** first runs something similar to the mini decompiler, then modifies the suggested tactics to produce a more natural proof script (11). For example, it cancels out sequences of `intros` and `revert`, inserts semicolons, and removes extra arguments to `apply` and `rewrite`. It can also take tactics from the proof engineer (like part of the old proof script) as hints, then iteratively replace tactics with those hints, checking for correctness. This makes it possible for suggested scripts to include custom tactics and decision procedures.

**From Qtac to Ltac.** The mini decompiler assumes more predictable versions of `rewrite` and `induction` than those in Coq. **Decompile** includes additional logic to reason about these tactics (11). For example, Qtac assumes that there is only one `rewrite` direction. Ltac has two `rewrite` directions, and so the decompiler infers the direction from the motive.

Qtac also assumes that both tactics take the inductive motive explicitly, while in Coq, both tactics infer the motive automatically. Consequentially, Coq sometimes fails to infer the correct motive. To handle induction, the decompiler strategically uses `revert` to manipulate the goal so that Coq can better infer the motive. To handle rewrites, it uses `simpl` to refold the goal before rewriting. Neither of these approaches is guaranteed to work, so the proof engineer may sometimes need to

tweak the suggested proof script appropriately. Even if we pass Coq’s induction principle an explicit motive, Coq still sometimes fails due to unrepresented assumptions. Long term, using another tactic like `change` or `refine` before applying these tactics may help with cases for which Coq cannot infer the correct motive.

**From  $CIC_\omega$  to Coq.** Scaling the decompiler to Coq introduces `let` bindings, which are generated by tactics like `rewrite in`, `apply in`, and `pose`. **Decompile** implements (11) support for `rewrite in` and `apply in` similarly to how it supports `rewrite` and `apply`, except that it ensures that the unmanipulated hypothesis does not occur in the body of the `let` expression, it swaps the direction of the `rewrite`, and it recurses into any generated subgoals. In all other cases, it uses `pose`, a catch-all for `let` bindings.

**Forfeiting Soundness.** While there is a way to always produce a correct proof script, **Decompile** deliberately forfeits soundness to suggest more useful tactics. For example, it may suggest the `induction` tactic, but leave the step of motive inference to the proof engineer. We have found these suggested tactics easier to work with (Section 4.6). Note that in the case the suggested proof script is not quite correct, it is still possible to use the generated proof term directly.

**Pretty Printing.** After decompiling proof terms, **Decompile** pretty prints the result (11). Like the mini decompiler, **Decompile** represents its output using a predefined grammar of Ltac tactics, albeit one that is larger than Qtac, and that also includes tacticals. It maintains the recursive proof structure for formatting. PUMPKIN Pi keeps all output terms from **Transform** in the Coq environment in case the decompiler does not succeed. Once the proof engineer has the new proof, she can remove the old one.

## 4.6 RESULTS

This section summarizes eight case studies using PUMPKIN Pi, corresponding to the eight rows in Table 1. These case studies highlight PUMPKIN Pi’s flexibility in handling diverse scenarios, the success of automatic configuration for better workflow integration, the preliminary success of the prototype decompiler, and clear paths to better serving proof engineers. Detailed walkthroughs are in the code.

**Algebraic Ornaments: Lists to Packed Vectors.** The transformation in PUMPKIN Pi is a generalization of the transformation from DEVoid. DEVoid supported proof reuse across *algebraic ornaments*, which describe relations between two inductive types, where one type is the other indexed by a fold [56]. A standard example is the relation between a list and a length-indexed vector (Figure 11).

PUMPKIN Pi implements a search procedure for automatic configuration of algebraic ornaments. The result is all functionality from

Class	Config.	Examples	Sav.	Repair Tools
Algebraic Ornaments	Auto	List to Packed Vector, hs-to-coq ③	☺	PUMPKIN Pi, DEVOID, UP
		List to Packed Vector, Std. Library ①⑥	☺	PUMPKIN Pi, DEVOID, UP
Unpack Sigma Types	Auto	Vector of Particular Length, hs-to-coq ③	☺	PUMPKIN Pi, UP
Tuples & Records	Auto	Simple Records ⑬	☺	PUMPKIN Pi, UP
		Parameterized Records ⑰	☺	PUMPKIN Pi, UP
		Industrial Use ⑱	☺	PUMPKIN Pi, UP
Permute Constructors	Auto	List, Standard Library ①	☺	PUMPKIN Pi, UP
		Modifying a PL, REPLICA Benchmark ①	☺	PUMPKIN Pi, UP
		Large Ambiguous Enum ①	☺	PUMPKIN Pi, UP
Add new Constructors	Mixed	PL Extension, REPLICA Benchmark ⑲	☺	PUMPKIN Pi
Factor out Constructors	Manual	External Example ②	☺	PUMPKIN Pi, UP
Permute Hypotheses	Manual	External Example ⑳	☺	PUMPKIN Pi, UP
Change Ind. Structure	Manual	Unary to Binary, Classic Benchmark ⑤	☺	PUMPKIN Pi, Magaud
		Vector to Finite Set, External Example ㉑	☺	PUMPKIN Pi

Table 1: Some changes using PUMPKIN Pi (left to right): class of changes, kind of configuration, examples, whether using PUMPKIN Pi saved development time relative to reference manual repairs (☺ if yes, ☺ if comparable, ☺ if no), and Coq tools we know of that support repair along (Repair) or automatic proof of (Search) the equivalence corresponding to each example. Tools considered are DEVOID [73], the Univalent Parametricity (UP) white-box transformation [78], and the tool from Magaud & Bertot 2000 [54]. PUMPKIN Pi is the only one that suggests tactics. More nuanced comparisons to these and more are in Section ??.

<pre> Inductive Term : Set :=   Var : Identifier → Term   Int : Z → Term   Eq : Term → Term → Term   Plus : Term → Term → Term   Times : Term → Term → Term   Minus : Term → Term → Term   Choose : Identifier → Term   → Term. </pre>	<pre> Inductive Term : Set :=   Var : Identifier → Term   Bool : Identifier → Term   Eq : Term → Term → Term   Int : Z → Term   Plus : Term → Term → Term   Times : Term → Term → Term   Minus : Term → Term → Term   Choose : Identifier → Term   → Term. </pre>
--	---

Figure 19: A simple language (left) and the same language with two swapped constructors and an added constructor (right).

DEVOID, plus tactic suggestions. In file ③, we used this to port functions and a proof from lists to vectors of *some* length, since  $\text{list } T \simeq \text{packed\_vect } T$ . The decompiler helped us write proofs in the order of hours that we had found too hard to write by hand, though the suggested tactics did need massaging.

**Unpack Sigma Types: Vectors of Particular Lengths.** In the same file ③, we then ported functions and proofs to vectors of a *particular* length, like  $\text{vector } T \ n$ . DEVOID had left this step to the proof engineer. We supported this in PUMPKIN Pi by chaining the previous change with an automatic configuration for unpacking sigma types. By composition, this transported proofs across the equivalence from Section ??.

Two tricks helped with workflow integration for this change: 1) have the search procedure view  $\text{vector } T \ n$  as  $\Sigma(v : \text{vector } T \ m). n = m$  for some  $m$ , then let PUMPKIN Pi instantiate those equalities via unification heuristics, and 2) generate a custom eliminator for combining list terms with length invariants. The resulting workflow works not just for lists and vectors, but for any algebraic ornament, automating manual effort from DEVOID. The suggested tactics were helpful for writing proofs in the order of hours that we had struggled with manually over the course of days, but only after massaging. More effort is needed to improve tactic suggestions for dependent types.

**Tuples & Records: Industrial Use.** An industrial proof engineer at the company Galois has been using PUMPKIN Pi in proving correct an implementation of the TLS handshake protocol. Galois had been using a custom solver-aided verification language to prove correct C programs, but had found that at times, the constraint solvers got stuck. They had built a compiler that translates their language into Coq’s specification language Gallina, that way proof engineers could finish stuck proofs interactively using Coq. However, due to language differences, they had found the generated Gallina programs and specifications difficult to work with.

The proof engineer used PUMPKIN Pi to port the automatically generated functions and specifications to more human-readable functions

and specifications, wrote Coq proofs about those functions and specifications, then used PUMPKIN Pi to port those proofs back to proofs about the original functions and specifications. So far, they have used at least three automatic configurations, but they most often used an automatic configuration for porting compiler-produced anonymous tuples to named records, as in file (18). The workflow was a bit nonstandard, so there was little need for tactic suggestions. The proof engineer reported an initial time investment learning how to use PUMPKIN Pi, followed by later returns.

**Permute Constructors: Modifying a Language.** The swapping example from Section 4.1 was inspired by benchmarks from the REPLICIA user study of proof engineers [71]. A change from one of the benchmarks is in Figure 19. The proof engineer had a simple language represented by an inductive type `Term`, as well as some definitions and proofs about the language. The proof engineer swapped two constructors in the language, and added a new constructor `Bool`.

This case study and the next case study break this change into two parts. In the first part, we used PUMPKIN Pi with automatic configuration to repair functions and proofs about the language after swapping the constructors (1). With a bit of human guidance to choose the permutation from a list of suggestions, PUMPKIN Pi repaired everything, though the original tactics would have also worked, so there was not a difference in development time.

**Add new Constructors: Extending a Language.** We then used PUMPKIN Pi to repair functions after adding the new constructor in Figure 19, separating out the proof obligations for the new constructor from the old terms (19). This change combined manual and automatic configuration. We defined an inductive type `Diff` and (using partial automation) a configuration to port the terms across the equivalence  $\text{Old.Term} + \text{Diff} \simeq \text{New.Term}$ . This resulted in case explosion, but was formulaic, and pointed to a clear path for automation of this class of changes. The repaired functions guaranteed preservation of the behavior of the original functions.

Adding constructors was less simple than swapping. For example, PUMPKIN Pi did not yet save us time over the proof engineer from the user study; fully automating the configuration would have helped significantly. In addition, the repaired terms were (unlike in the swap case) inefficient compared to human-written terms. For now, they make good regression tests for the human-written terms—in the future, we hope to automate the discovery of the more efficient terms, or use the refinement framework CoqEAL [20] to get between proofs of the inefficient and efficient terms.

**Factor out Constructors: External Example.** The change from Figure 10 came at the request of a non-author. We supported this using a manual configuration that described which constructor to map to true

and which constructor to map to `false` ②. The configuration was very simple for us to write, and the repaired tactics were immediately useful. The development time savings were on the order of minutes for a small proof development. Since most of the modest development time went into writing the configuration, we expect time savings would increase for a larger development.

**Permute Hypotheses: External Example.** The change in ②① came at the request of a different non-author (a cubical type theory expert), and shows how to use PUMPKIN Pi to swap two hypotheses of a type, since  $T1 \rightarrow T2 \rightarrow T3 \simeq T2 \rightarrow T1 \rightarrow T3$ . This configuration was manual. Since neither type was inductive, this change used the generic construction for any equivalence. This worked well, but necessitated some manual annotation due to the lack of custom unification heuristics for manual configuration, and so did not yet save development time, and likely still would not have had the proof development been larger. Supporting custom unification heuristics would improve this workflow.

**Change Inductive Structure: Unary to Binary.** In ⑤, we used PUMPKIN Pi to support a classic example of changing inductive structure: updating unary to binary numbers, as in Figure ?? . Binary numbers allow for a fast addition function, found in the Coq standard library. In the style of Magaud & Bertot 2000 [54], we used PUMPKIN Pi to derive a slow binary addition function that does not refer to `nat`, and to port proofs from unary to slow binary addition. We then showed that the ported theorems hold over fast binary addition.

The configuration for `N` used definitions from the Coq standard library for `DepConstr` and `DepElim` that had the desired behavior with no changes. `Iota` over the successor case was a rewrite by a lemma from the standard library that reduced the successor case of the eliminator that we used for `DepElim`:

```
N.peano_rect_succ : ∀ P p0 pS n,
  N.peano_rect P p0 pS (N.succ n) =
  pS n (N.peano_rect P p0 pS n).
```

The need for nontrivial `Iota` comes from the fact that `N` and `nat` have different inductive structures. By writing a manual configuration with this `Iota`, it was possible for us to implement this transformation that had been its own tool.

While porting addition from `nat` to `N` was automatic after configuring PUMPKIN Pi, porting proofs about addition took more work. Due to the lack of unification heuristics for manual configuration, we had to annotate the proof term to tell PUMPKIN Pi that implicit casts in the inductive cases of proofs were applications of `Iota` over `nat`. These annotations were formulaic, but tricky to write. Unification heuristics would go a long way toward improving the workflow.

After annotating, we obtained automatically repaired proofs about slow binary addition, which we found simple to port to fast binary



addition. We hope to automate this last step in the future using Co-qEAL. Repaired tactics were partially useful, but failed to understand custom eliminators like `N.peano_rect`, and to generate useful tactics for applications of `Iota`; both of these are clear paths to more useful tactics. The development time for this proof with PUMPKIN Pi was comparable to reference manual repairs by external proof engineers. Custom unification heuristics would help bring returns on investment for experts in this use case.

#### 4.7 CONCLUSION

Rehashing thesis and how we do it

What we got here beyond what we had in PUMPKIN PATCH, segue into next chapter



# 5

---

## RELATED WORK

---

### 5.1 PROGRAMS

#### *Program Refactoring*

Refactoring [60].

#### *Program Repair*

Adapting proofs to changes is essentially program repair for dependently typed languages. Program repair tools for languages with non-dependent type systems [67, 52, 49, 58, 64] may have applications in the context of a dependently typed language. Similarly, our work may have applications within program repair in these languages: Future applications of our approach may repurpose it to repair programs for functional languages.

#### *Ornaments*

DEVoid automates discovery of and lifting across algebraic ornaments in a higher-order dependently typed language. In the decade since the discovery of ornaments [56], there have been a number of formalizations and embedded implementations of ornaments [24, 45, 25, 46, 23]. DEVoid is the first tool for ornamentation to operate over a non-embedded dependently typed language. It essentially moves the automation-heavy approach of Ornamentation in ML [81], which operates on non-embedded ML code, into the type theory that forms the basis of theorem provers like Coq. In doing so, it takes advantage of the properties of algebraic ornaments [56]. It also introduces the first search algorithm to identify ornaments, which in the past was identified as a “gap” in the literature [46].

#### *Programming by Example*

Our approach generalizes an example that the programmer provides. This is similar to programming by example, a subfield of program

synthesis [38]. This field addresses different challenges in different logics, but may drive solutions to similar problems in a dependently typed language.

### *Differencing & Incremental Computation*

Existing work in differencing and incremental computation may help improve our semantic differencing component. Type-directed differencing [63] finds differences in algebraic data types. Semantics-based change impact analysis [5] models semantic differences between documents. Differential assertion checking [48] analyzes different versions of a program for relative correctness with respect to a specification. Incremental  $\lambda$ -calculus [13] introduces a general model for program changes. All of these may be useful for improving semantic differencing.

## 5.2 PROOFS

### *Proof Reuse*

Our approach reimagines the problem of proof reuse in the context of proof automation. While we focus on changes that occur over time, traditional proof reuse techniques can help improve our approach. Existing work in proof reuse focuses on transferring proofs between isomorphisms, either through extending the type system [8] or through an automatic method [55]. This is later generalized and implemented in Isabelle [40] and Coq [83, 76]; later methods can also handle implications. Integrating a transfer tactic with a proof patch finding tool will create an end-to-end tool that can both find patches and apply them automatically.

Proof reuse for extended inductive types [10] adapts proof obligations to structural changes in inductive types. Later work [65] proposes a method to generate proofs for new constructors. These approaches may be useful when extending the differencing component to handle structural changes. Existing work in theorem reuse and proof generalization [33, 68, 43] abstracts existing proofs for reusability, and may be useful for improving the abstraction component. Our work focuses on the components critical to searching for patches; these complementary approaches can drive improvements to the components.

DEVOID identifies and lifts proofs along a specific equivalence similar to that from existing ornaments work [46]. The need to automatically lift functions and proofs across equivalences and other relations is a long-standing challenge for proof engineers [54, 9, 53, 41, 84, 21]. The univalence axiom from Homotopy Type Theory [80] enables transparent transport of proofs; cubical type theory [19] gives univalence a constructive interpretation.

Similarly, our work is related to CoqEAL [21], which transfers functions along arbitrary relations between types. As these relations do not necessarily need to be equivalences, this framework is more general than our work. Similar tradeoffs between automation and generality apply: CoqEAL produces functions that refer to the old type, and does not yet support automatic inference of relations. In addition, CoqEAL currently only supports automatic transfer of functions, and does not yet handle proofs.

These tools may provide an alternative backend for DEVoid. Furthermore, our search algorithm may help discover relations that make these tools easier to use, and our lifting algorithm may help improve automation and efficiency for certain relations in these tools.

The problem that we solve is fundamentally about proof reuse, which applies software reuse principles to ITPs. There is a wealth of work in proof reuse, from tactic languages [34] and logical frameworks [14], to tools for proof abstraction and generalization [69, 44], to domain-specific methodologies [26] and frameworks [27].

DEVoid focuses on the specific problem of reuse when adding fully-determined indices to types. Other approaches to this problem include combinators which definitionally reduce to desirable terms [32] in the language Cedille, and automatic generation of conversion functions in Ghostbuster [57] for GADTs in Haskell. Our work focuses on a type theory different from both of these, in which the properties that allow for such combinators in Cedille are not present, and in which dependent types introduce challenges not present in Haskell.

DEVoid is not the first tool to combine search with reuse. Optician [62] synthesizes bidirectional string transformations; a similar approach may help extend tooling to handle transformations for low-level data. PUMPKIN PATCH [72] searches the difference in proofs for patches that can be used to repair proofs broken by changes; DEVoid uses a similar approach to identify functions that form an equivalence. The resulting tools are complementary: DEVoid supports the addition of indices and hypotheses, which PUMPKIN PATCH does not support; PUMPKIN PATCH supports changes in values, which DEVoid does not support.

### *Proof Evolution*

There is a small body of work on change and dependency management for verification, both to evaluate impact of potential changes and maximize reuse [42, 4] and to optimize build performance [15]. These approaches may help isolate changes, which is necessary to identify future benchmarks, integrate with CI systems, and fully support version updates.

*Proof Refactoring**Proof Repair**Proof Design*

Existing proof engineering work addresses brittleness by planning for changes [82] and designing theorems and proofs that make maintenance less of an issue. Design principles for specific domains (such as formal metatheory [6, 28, 29]) can make verification more tractable. CertiKOS [37] introduces the idea of a deep specification to ease verification of large systems. These design principles and frameworks are complementary to our approach. Even when programmers use informed design principles, changes outside of the programmer’s control can break proofs; our approach addresses these changes.

*Proof Automation*

We address a missed opportunity in proof automation for ITP: searching for patches that can fix broken proofs. This is complementary to existing automation techniques. Nonetheless, there is a wealth of work in proof automation that makes proofs more resilient to change. Powerful tactics like *crush* [17] can make proofs more resilient to changes. Hammers like Isabelle’s *sledgehammer* [66] can make proofs agnostic to some low-level changes. Recent work [22] paves the way for a hammer in Coq. Even the most powerful tactics cannot address all changes; our hope is to open more possibilities for automation.

Powerful project-specific tactics [17, 16] can help prevent low-level maintenance tasks. Writing these tactics requires good engineering [36] and domain-specific knowledge, and these tactics still sometimes break in the face of change. A future patching tool may be able to repair tactics; the debugging process for adapting a tactic is not too dissimilar to providing an example to a tool.

Rippling [12] is a technique for automating inductive proofs that uses restricted rewrite rules to guide the inductive hypothesis toward the conclusion; this may guide improvements to the differencing, abstraction, and specialization components. The abstraction and factoring components address specific classes of unification problems; recent developments to higher-order unification [61] may help improve these components. Lean [75] introduces the first congruence closure algorithm for dependent type theory that relies only on the Uniqueness of Identity Proofs (UIP) axiom. While UIP is not fundamental to Coq, it is frequently assumed as an axiom; when it is, it may be tractable to use a similar algorithm to improve the tool.

GALILEO [11] repairs faulty physics theories in the context of a classical higher-order logic (HOL); there is preliminary work extend-

ing this style of repair to mathematical proofs. Knowledge-sharing methods [35] can adapt some proofs across different representations of HOL. These complementary approaches may guide extensions to support decidable domains and classical logics.





# 6

---

## CONCLUSIONS & FUTURE WORK

---

Reflect on thesis statement and explain how we got it exactly now that you know everything

But I want to spend the rest of this thesis talking about the next era of verification so I can write out a bunch of ideas for students who might want to work with me

### THE NEXT ERA: PROOF ENGINEERING FOR ALL

Future Work from many papers, plus research statement, DARPA thoughts, plus more, but trimmed down a lot

What I want in the long run, how this all fits in, is a world of proof engineering for all. From research statement, three rings (four including experts in the center).

And what we have so far with my thesis is a world where it's easier for experts and a bit easier for practitioners, but there's still a lot left to go building on it.

So here are 12 short future project summaries that reach each of these tiers, building that world. Super please contact me if any of these seem fun to you.

#### *Proof Engineering for Experts*

Unifying theme: lateral reach. Some examples:

**MORE PROOF ASSISTANTS** Thoughts from PUMPKIN Pi on Isabelle/HOL, future work from PUMPKIN PATCH.

**MORE CHANGES** Version updates, isolating large changes (PUMPKIN PATCH), relations more general than equivalences (PUMPKIN Pi).

**MORE STYLES** ML for decompiler (PUMPKIN Pi, REPLICA) : more for diverse proof styles (PUMPKIN PATCH). Note that this is a WIP, but sketch out project, challenges, future ideas, expectations, evaluation a bit.

*Proof Engineering for Practitioners*

Unifying theme: usability. Some examples:

**AUTOMATION** More search procedures for automatic configuration, e-graphs from PUMPKIN Pi, custom unification heuristics.

**INTEGRATION** IDE & CI integration, HCI for repair.

**EVALUATION** repair challenge, user studies ideas (PUMPKIN PATCH, REPLICA, panel w/ Benjamin Pierce, QED at large). (maybe look for more ideas, this can be merged with integration if need be).

*Proof Engineering for Software Engineers*

Unifying theme: mixed methods verification, or the 2030 vision from Twitter thread. Some examples:

**GRADUAL VERIFICATION** A continuum from testing to verification, tools to help with that.

**TOOL-ASSISTED PROOF DEVELOPMENT** Tool-assisted development to follow good design principles for verification (James Wilcox conversation, final REPLICA takeaway).

**SPECIFICATION INFERENCE** Analysis to infer specs (TA1).

*Proof Engineering for New Domains*

Unifying theme: collaboration, new abstractions for new domains). Some examples:

**MACHINE LEARNING** Fairification & other ML correctness properties. Some stuff here but more.

**CRYPTOGRAPHY** Lots of stuff here but not thinking broadly enough. What about cryptographic proof systems? ZK and beyond. Recall email thread.

**SOMETHING ELSE** Look for more in survey paper, email, DARPA TAs, Twitter. Healthcare perhaps?

---

## BIBLIOGRAPHY

---

- [1] Coq reference manual, section 8.9: Controlling automation, 2017.
- [2] User A. Software foundations solution, 2017.
- [3] Carlo Angiuli, Evan Cavallo, Anders Mörtberg, and Max Zeuner. Internalizing representation independence with univalence, 2020.
- [4] Serge Autexier, Dieter Hutter, and Till Mossakowski. Verification, induction termination analysis. chapter Change Management for Heterogeneous Development Graphs, pages 54–80. Springer-Verlag, Berlin, Heidelberg, 2010.
- [5] Serge Autexier and Normen Müller. Semantics-based change impact analysis for heterogeneous collections of documents. In *Proceedings of the 10th ACM Symposium on Document Engineering, DocEng '10*, pages 97–106, New York, NY, USA, 2010. ACM.
- [6] Brian Aydemir, Arthur Charguéraud, Benjamin C. Pierce, Randy Pollack, and Stephanie Weirich. Engineering formal metatheory. In *Proceedings of the 35th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '08*, pages 3–15, New York, NY, USA, 2008. ACM.
- [7] User B. Software foundations solution, 2017.
- [8] Gilles Barthe and Olivier Pons. Type isomorphisms and proof reuse in dependent type theory. In *Proceedings of the 4th International Conference on Foundations of Software Science and Computation Structures, FoSSaCS '01*, pages 57–71, London, UK, UK, 2001. Springer-Verlag.
- [9] Gilles Barthe and Olivier Pons. Type isomorphisms and proof reuse in dependent type theory. In *International Conference on Foundations of Software Science and Computation Structures*, pages 57–71. Springer, 2001.
- [10] Olivier Boite. Proof reuse with extended inductive types. In Konrad Slind, Annette Bunker, and Ganesh Gopalakrishnan, editors, *Theorem Proving in Higher Order Logics: 17th International Conference, TPHOLs 2004, Park City, Utah, USA, September 14-17, 2004. Proceedings*, pages 50–65, Berlin, Heidelberg, 2004. Springer Berlin Heidelberg.
- [11] Alan Bundy. The interaction of representation and reasoning. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 469(2157), 2013.

- [12] Alan Bundy, David Basin, Dieter Hutter, and Andrew Ireland. *Rippling: Meta-Level Guidance for Mathematical Reasoning*. Cambridge University Press, New York, NY, USA, 2005.
- [13] Yufei Cai, Paolo G. Giarrusso, Tillmann Rendel, and Klaus Ostermann. A theory of changes for higher-order languages: Incrementalizing  $\lambda$ -calculi by static differentiation. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '14*, pages 145–155, New York, NY, USA, 2014. ACM.
- [14] Joshua E. Caplan and Mehdi T. Harandi. A logical framework for software proof reuse. In *ACM SIGSOFT Software Engineering Notes*, volume 20, pages 106–113. ACM, 1995.
- [15] Ahmet Celik, Karl Palmskog, and Milos Gligoric. icoq: Regression proof selection for large-scale verification projects. In *Proceedings of the 32Nd IEEE/ACM International Conference on Automated Software Engineering, ASE 2017*, pages 171–182, Piscataway, NJ, USA, 2017. IEEE Press.
- [16] Adam Chlipala. The Bedrock structured programming system: Combining generative metaprogramming and Hoare logic in an extensible program verifier. In *Proceedings of the 18th ACM SIGPLAN International Conference on Functional Programming, ICFP '13*, pages 391–402, New York, NY, USA, 2013. ACM.
- [17] Adam Chlipala. *Certified Programming with Dependent Types - A Pragmatic Introduction to the Coq Proof Assistant*. MIT Press, 2013.
- [18] Adam Chlipala. Library equality, 2017.
- [19] Cyril Cohen, Thierry Coquand, Simon Huber, and Anders Mörtberg. Cubical type theory: a constructive interpretation of the univalence axiom. *arXiv preprint arXiv:1611.02108*, 2016.
- [20] Cyril Cohen, Maxime Dénès, and Anders Mörtberg. Refinements for Free! In *Certified Programs and Proofs*, pages 147 – 162, Melbourne, Australia, December 2013.
- [21] Cyril Cohen and Damien Rouhling. A refinement-based approach to large scale reflection for algebra. In *JFLA 2017 - Vingt-huitième Journées Francophones des Langages Applicatifs*, Gourette, France, January 2017.
- [22] Łukasz Czajka and Cezary Kaliszyk. Hammer for coq: Automation for dependent type theory. *Journal of Automated Reasoning*, 61(1):423–453, Jun 2018.
- [23] Pierre-Évariste Dagand. The essence of ornaments. *Journal of Functional Programming*, 27, 2017.

- [24] Pierre-Evariste Dagand and Conor McBride. A categorical treatment of ornaments. In *Proceedings of the 2013 28th Annual ACM/IEEE Symposium on Logic in Computer Science, LICS '13*, pages 530–539, Washington, DC, USA, 2013. IEEE Computer Society.
- [25] Pierre-Evariste Dagand and Conor McBride. Transporting functions across ornaments. *Journal of functional programming*, 24(2-3):316–383, 2014.
- [26] Benjamin Delaware, William Cook, and Don Batory. Product lines of theorems. In *Proceedings of the 2011 ACM International Conference on Object Oriented Programming Systems Languages and Applications, OOPSLA '11*, pages 595–608, New York, NY, USA, 2011. ACM.
- [27] Benjamin Delaware, Bruno C. d. S. Oliveira, and Tom Schrijvers. Meta-theory à la carte. In *Proceedings of the 40th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '13*, pages 207–218, New York, NY, USA, 2013. ACM.
- [28] Benjamin Delaware, Bruno C. d. S. Oliveira, and Tom Schrijvers. Meta-theory à la carte. In *Proceedings of the 40th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '13*, pages 207–218, New York, NY, USA, 2013. ACM.
- [29] Benjamin Delaware, Steven Keuchel, Tom Schrijvers, and Bruno C.d.S. Oliveira. Modular monadic meta-theory. In *Proceedings of the 18th ACM SIGPLAN International Conference on Functional Programming, ICFP '13*, pages 319–330, New York, NY, USA, 2013. ACM.
- [30] Richard A. DeMillo, Richard J. Lipton, and Alan J. Perlis. Social processes and proofs of theorems and programs. In *Proceedings of the 4th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages, POPL '77*, pages 206–214, New York, NY, USA, 1977. ACM.
- [31] Maxime Dénes. Coq 8.7 beta 1 is out, 2017.
- [32] Larry Diehl, Denis Firsov, and Aaron Stump. Generic zero-cost reuse for dependent types. *CoRR*, abs/1803.08150, 2018.
- [33] Amy Felty and Douglas Howe. Generalization and reuse of tactic proofs. In Frank Pfenning, editor, *Logic Programming and Automated Reasoning: 5th International Conference, LPAR '94*, pages 1–15, Berlin, Heidelberg, 1994. Springer Berlin Heidelberg.

- [34] Amy Felty and Douglas Howe. Generalization and reuse of tactic proofs. In *International Conference on Logic for Programming Artificial Intelligence and Reasoning*, pages 1–15. Springer, 1994.
- [35] Thibault Gauthier and Cezary Kaliszyk. Matching concepts across HOL libraries. In Stephen Watt, James Davenport, Alan Sexton, Petr Sojka, and Josef Urban, editors, *CICM '14*, volume 8543 of *LNCS*, pages 267–281. Springer Verlag, 2014.
- [36] Georges Gonthier, Beta Ziliani, Aleksandar Nanevski, and Derek Dreyer. How to make ad hoc proof automation less ad hoc. In *Proceedings of the 16th ACM SIGPLAN International Conference on Functional Programming, ICFP '11*, pages 163–175, New York, NY, USA, 2011. ACM.
- [37] Ronghui Gu, Zhong Shao, Hao Chen, Xiongnan (Newman) Wu, Jieung Kim, Vilhelm Sjöberg, and David Costanzo. Certikos: An extensible architecture for building certified concurrent OS kernels. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pages 653–669, GA, 2016. USENIX Association.
- [38] Sumit Gulwani, Oleksandr Polozov, and Rishabh Singh. Program synthesis. *Foundations and Trends in Programming Languages*, 4(1-2):1–119, 2017.
- [39] Martin Hofmann. Syntax and semantics of dependent types. In *Semantics and Logics of Computation*, pages 79–130. Cambridge University Press, 1997.
- [40] Brian Huffman and Ondřej Kunčar. Lifting and transfer: A modular design for quotients in Isabelle/HOL. In *Certified Programs and Proofs: Third International Conference, CPP 2013*, pages 131–146, Cham, 2013. Springer International Publishing.
- [41] Brian Huffman and Ondřej Kunčar. Lifting and Transfer: A modular design for quotients in Isabelle/HOL. In *International Conference on Certified Programs and Proofs*, pages 131–146. Springer, 2013.
- [42] D. Hutter. Management of change in structured verification. In *ASE 2000*, pages 23–31, Sept 2000.
- [43] Einar Broch Johnsen and Christoph Lüth. Theorem reuse by proof term transformation. In *Theorem Proving in Higher Order Logics: 17th International Conference, TPHOLs 2004, Park City, Utah, USA, September 14-17, 2004. Proceedings*, pages 152–167. Springer, Berlin, Heidelberg, 2004.

- [44] Einar Broch Johnsen and Christoph Lüth. Theorem reuse by proof term transformation. In *International Conference on Theorem Proving in Higher Order Logics*, pages 152–167. Springer, 2004.
- [45] Hsiang-Shang Ko and Jeremy Gibbons. Relational algebraic ornaments. In *Proceedings of the 2013 ACM SIGPLAN workshop on Dependently-typed programming*, pages 37–48. ACM, 2013.
- [46] Hsiang-Shang Ko and Jeremy Gibbons. Programming with ornaments. *Journal of Functional Programming*, 27, 2016.
- [47] Matej Kosik. Coq pull request # 652: Put all plugins behind an “api”, 2017.
- [48] Shuvendu Lahiri, Kenneth McMillan, , and Chris Hawblitzel. Differential assertion checking. In *Foundations of Software Engineering (FSE’13)*. ACM, August 2013.
- [49] Xuan-Bach D. Le, Duc-Hiep Chu, David Lo, Claire Le Goues, and Willem Visser. S3: Syntax- and semantic-guided repair synthesis via programming by examples. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering, ESEC/FSE 2017*, pages 593–604, New York, NY, USA, 2017. ACM.
- [50] Xavier Leroy. Commit to compcert: lib/integers.v, 2013.
- [51] letouzey. Commit to coq: change definition of divide (compat with znumtheory), 2011.
- [52] Fan Long and Martin Rinard. Automatic patch generation by learning correct code. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL ’16*, pages 298–312, New York, NY, USA, 2016. ACM.
- [53] Nicolas Magaud. Changing data representation within the Coq system. In *International Conference on Theorem Proving in Higher Order Logics*, pages 87–102. Springer, 2003.
- [54] Nicolas Magaud and Yves Bertot. Changing data structures in type theory: A study of natural numbers. In *International Workshop on Types for Proofs and Programs*, pages 181–196. Springer, 2000.
- [55] Nicolas Magaud and Yves Bertot. Changing data structures in type theory: A study of natural numbers. In Paul Callaghan, Zhaohui Luo, James McKinna, and Robert Pollack, editors, *Types for Proofs and Programs: International Workshop, TYPES 2000*, pages 181–196, Berlin, Heidelberg, 2002. Springer Berlin Heidelberg.
- [56] Conor McBride. Ornamental algebras, algebraic ornaments, 2011.

- [57] Trevor L. McDonell, Timothy A. K. Zakian, Matteo Cimini, and Ryan R. Newton. Ghostbuster: A tool for simplifying and converting GADTs. In *Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming, ICFP 2016*, pages 338–350, New York, NY, USA, 2016. ACM.
- [58] Sergey Mechtaev, Jooyong Yi, and Abhik Roychoudhury. Angelix: Scalable multiline program patch synthesis via symbolic analysis. In *Proceedings of the 38th International Conference on Software Engineering, ICSE '16*, pages 691–701, New York, NY, USA, 2016. ACM.
- [59] Guillaume Melquiond. Commit to coq: Make izr use a compact representation of integers, 2017.
- [60] Tom Mens and Tom Tourwé. A survey of software refactoring. *IEEE Trans. Softw. Eng.*, 30(2):126–139, February 2004.
- [61] Dale Miller and Gopalan Nadathur. *Programming with Higher-Order Logic*. Cambridge University Press, New York, NY, USA, 1st edition, 2012.
- [62] Anders Miltner, Kathleen Fisher, Benjamin C Pierce, David Walker, and Steve Zdancewic. Synthesizing bijective lenses. *Proceedings of the ACM on Programming Languages*, 2(POPL):1, 2017.
- [63] Victor Cacciari Miraldo, Pierre-Évariste Dagand, and Wouter Swierstra. Type-directed diffing of structured data. In *Proceedings of the 2Nd ACM SIGPLAN International Workshop on Type-Driven Development, TyDe 2017*, pages 2–15, New York, NY, USA, 2017. ACM.
- [64] Martin Monperrus. Automatic Software Repair: a Bibliography. *ACM Computing Surveys*, 2017.
- [65] Anne Mulhern. Proof weaving. In *In Proceedings of the First Informal ACM SIGPLAN Workshop on Mechanizing Metatheory*, 2006.
- [66] Lawrence C. Paulson and Jasmin Christian Blanchette. Three years of experience with Sledgehammer, a practical link between automatic and interactive theorem provers. In G. Sutcliffe, S. Schulz, and E. Ternovska, editors, *International Workshop on the Implementation of Logics (IWIL 2010)*, volume 2 of *EPiC Series*, pages 1–11. EasyChair, 2012.
- [67] Yu Pei, Carlo A. Furia, Martin Nordio, and Bertrand Meyer. Automatic program repair by fixing contracts. In *Proceedings of the 17th International Conference on Fundamental Approaches to Software Engineering - Volume 8411*, pages 246–260, New York, NY, USA, 2014. Springer-Verlag New York, Inc.



- [68] Olivier Pons. Generalization in type theory based proof assistants. *TYPES '00*, pages 217–232, 12 2000.
- [69] Olivier Pons. Generalization in type theory based proof assistants. In *International Workshop on Types for Proofs and Programs*, pages 217–232. Springer, 2000.
- [70] Talia Ringer, Karl Palmskog, Ilya Sergey, Milos Gligoric, and Zachary Tatlock. Qed at large: A survey of engineering of formally verified software. *Foundations and Trends® in Programming Languages*, 5(2-3):102–281, 2019.
- [71] Talia Ringer, Alex Sanchez-Stern, Dan Grossman, and Sorin Lerner. Replica: Repl instrumentation for coq analysis. In *Proceedings of the 9th ACM SIGPLAN International Conference on Certified Programs and Proofs*, CPP 2020, pages 99–113, 01 2020.
- [72] Talia Ringer, Nathaniel Yazdani, John Leo, and Dan Grossman. Adapting proof automation to adapt proofs. In *Proceedings of the 7th ACM SIGPLAN International Conference on Certified Programs and Proofs*, pages 115–129. ACM, 2018.
- [73] Talia Ringer, Nathaniel Yazdani, John Leo, and Dan Grossman. Ornaments for Proof Reuse in Coq. In John Harrison, John O’Leary, and Andrew Tolmach, editors, *10th International Conference on Interactive Theorem Proving (ITP 2019)*, volume 141 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 26:1–26:19, Dagstuhl, Germany, 2019. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.
- [74] Valentin Robert. *Front-end tooling for building and maintaining dependently-typed functional programs*. PhD thesis, UC San Diego, 2018.
- [75] Daniel Selsam and Leonardo de Moura. Congruence closure in intensional type theory. In Nicola Olivetti and Ashish Tiwari, editors, *Automated Reasoning: 8th International Joint Conference, IJCAR 2016*, pages 99–115, Cham, 2016. Springer International Publishing.
- [76] Nicolas Tabareau, Éric Tanter, and Matthieu Sozeau. Equivalences for Free! working paper or preprint, July 2017.
- [77] Nicolas Tabareau, Éric Tanter, and Matthieu Sozeau. Equivalences for free: Univalent parametricity for effective transport. *Proc. ACM Program. Lang.*, 2(ICFP):92:1–92:29, July 2018.
- [78] Nicolas Tabareau, Éric Tanter, and Matthieu Sozeau. The marriage of univalence and parametricity, 2019.

- [79] Amin Timany and Bart Jacobs. First steps towards cumulative inductive types in CIC. In *ICTAC*, 2015.
- [80] The Univalent Foundations Program. *Homotopy Type Theory: Univalent Foundations of Mathematics*. <https://homotopytypetheory.org/book>, Institute for Advanced Study, 2013.
- [81] Thomas Williams and Didier Rémy. A principled approach to ornamentation in ML. *Proc. ACM Program. Lang.*, 2(POPL):21:1–21:30, December 2017.
- [82] Doug Woos, James R. Wilcox, Steve Anton, Zachary Tatlock, Michael D. Ernst, and Thomas Anderson. Planning for change in a formal verification of the raft consensus protocol. In *Proceedings of the 5th ACM SIGPLAN Conference on Certified Programs and Proofs, CPP 2016*, pages 154–165, New York, NY, USA, 2016. ACM.
- [83] Théo Zimmermann and Hugo Herbelin. Automatic and transparent transfer of theorems along isomorphisms in the coq proof assistant. *CoRR*, abs/1505.05028, 2015.
- [84] Theo Zimmermann and Hugo Herbelin. Automatic and transparent transfer of theorems along isomorphisms in the Coq proof assistant. *arXiv preprint arXiv:1505.05028*, 2015.