

# Adapting Proof Automation to Adapt Proofs

Talia Ringer  
University of Washington, USA

John Leo  
Halfaya Research, USA

Nathaniel Yazdani  
University of Washington, USA

Dan Grossman  
University of Washington, USA

## Abstract

We extend proof automation in an interactive theorem prover to analyze *changes* in specifications and proofs. Our approach leverages the history of changes to specifications and proofs to search for a patch that can be applied to other specifications and proofs that need to change in analogous ways.

We identify and implement five core components that are key to searching for a patch. We build a patch finding procedure from these components, which we configure for various classes of changes. We implement this procedure in a Coq plugin as a proof-of-concept and use it on real Coq code to change specifications, port definitions of a type, and update the Coq standard library. We show how our findings help drive a future that moves the burden of dealing with the brittleness of small changes in an interactive theorem prover away from the programmer and into automated tooling.

## ACM Reference Format:

Talia Ringer, Nathaniel Yazdani, John Leo, and Dan Grossman. 2017. Adapting Proof Automation to Adapt Proofs. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 16 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 Introduction

Proof automation makes verification more accessible to programmers, but it is often intractable without programmer guidance. In *interactive theorem proving* (ITP), the programmer guides the proof search process. The guidance reduces the search space to make proof automation more tractable. This in turn helps the programmer, who does not have to manually verify the entire system.

Despite automation, programming in these proof assistants is brittle: Even a minor change to a definition or theorem can break many dependent proofs. This is a major source of development inefficiency in proof assistants based on dependent type theory [13, 29, 64].

Traditional proof automation does not consider how proofs, definitions, and theorems change over time. Instead, it is driven by the state of the current proof, sometimes with supplementary information from other

proofs, definitions, and theorems (as in hint databases [3] and rippling [20]). This puts the burden of dealing with brittleness on the programmer.

We present a new approach to proof automation that accounts for breaking changes. In our approach, the programmer guides proof search by providing an *example* of how to adapt proofs to changes in definitions or theorems. A tool then generalizes the example adaptation into a *reusable patch* that the programmer can use to fix other broken proofs.

In doing so, we chart a path for a future that moves the burden of brittleness away from the programmer and into proof automation. Programmers typically address brittleness through design principles that make proofs resilient to change [13, 29, 64], or through program-specific proof automation [24]. These techniques, while useful, have limitations: Planning a verification effort around future change is challenging, and program-specific proof automation requires specialized knowledge. The programmer's ability to anticipate likely changes determines the robustness of both techniques in the face of change. Even then, many breaking changes are outside of the programmer's control. Updating proof assistant versions, for example, can break proofs regardless of planning or automation [55].

We identify a set of core components that are critical to searching an example for patches in Coq. We use the components to build a procedure for finding patches in a Coq plugin as a proof-of-concept. Case studies on real projects like CompCert [43] and the Coq standard library suggest that proof patches are useful for realistic scenarios, and that it is simple to compose the core components to handle different classes of changes. We test the boundaries of the core components on a suite of tests and show how our findings can help build the future we envision.

In summary, we contribute the following:

1. We identify a set of core components that are key to searching for patches.
2. We demonstrate that these core components are useful on real changes in Coq code.
3. We explore how to drive a future that moves the burden of change in ITP away from the programmer and into automated tooling.

```

Definition IZR (z:Z) : R := match z with
| Z0 => 0
| Zpos n => INR (Pos.to_nat n)
| Zneg n => - INR (Pos.to_nat n)
end.

```

```

Definition IZR (z:Z) : R := match z with
| Z0 => 0
| Zpos n => IPR n
| Zneg n => - IPR n
end.

```

**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.

## 2 Automation, Reimagined

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.

Currently, the burden of change largely falls on programmers. 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 version and the new version for a *reusable patch* that can fix broken proofs.

**Status Quo** 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.

Consider a commit from the upcoming Coq 8.7 release [49]. 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 : ∀ 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 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.

**Our Vision** When a user updates Coq, a 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 Generating Reusable Patches

We identify five components that are key to finding reusable patches. We implement these components in a prototype Coq plugin, which we call PUMPKIN PATCH (Proof Updater Mechanically Passing Knowledge Into New Proofs, Assisting The Coq Hacker), or PUMPKIN.<sup>1</sup>

We focus the PUMPKIN prototype on the proof assistant Coq [2]. In Coq, each theorem is a type, and a proof of that theorem is a term that inhabits that type. Rather than write proof terms directly, users write proof scripts in a high-level tactic language; Coq then uses these scripts to guide search for a proof term.

The PUMPKIN repository contains a detailed user guide. To use PUMPKIN, the programmer modifies a single proof script to provide an *example* of how to adapt a proof to a change. PUMPKIN generalizes the example adaptation into a *reusable patch*: a function that can be used to fix other broken proofs, which PUMPKIN defines as a Coq term. We focus on the problem of *finding* these patches; we discuss how to extend PUMPKIN to automatically *apply* patches it finds in Section 7.

We motivate the core components (Section 3.1) and describe how they compose to find patches (Section 3.2). We find patches for real code in Section 4, and we describe our implementation in Section 5.

### 3.1 Motivating the Core

We identify and isolate five core components critical to finding reusable patches:

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

The semantic differencing component finds the difference between two terms, which produces a *candidate*

<sup>1</sup><http://github.com/uwplse/PUMPKIN-PATCH/tree/cpp18>

1 <b>Theorem</b> old: $\forall (n\ m\ p : \text{nat}), n \leq m \rightarrow m \leq p \rightarrow$	1 <b>Theorem</b> new: $\forall (n\ m\ p : \text{nat}), n \leq m \rightarrow m \leq p \rightarrow$
2 $n \leq p + 1.$ <span style="float: right;">(* P p *)</span>	2 $n \leq p.$ <span style="float: right;">(* P' p *)</span>
3 <b>Proof.</b>	3 <b>Proof.</b>
4 intros. induction H0.	4 intros. induction H0.
5 - auto with arith.	5 - auto with arith.
6 - constructor. auto.	6 - constructor. auto.
7 Qed.	7 Qed.
8	8
9 fun (n m p : nat) (H : n ≤ m) (H0 : m ≤ p) =>	9 fun (n m p : nat) (H : n ≤ m) (H0 : m ≤ p) =>
10 le_ind	10 le_ind
11 m <span style="float: right;">(* m *)</span>	11 m <span style="float: right;">(* m *)</span>
12 (fun p0 => n ≤ p0 + 1) <span style="float: right;">(* P *)</span>	12 (fun p0 => n ≤ p0) <span style="float: right;">(* P' *)</span>
13 (le_plus_trans n m 1 H) <span style="float: right;">(* : P m *)</span>	13 H <span style="float: right;">(* : P' m *)</span>
14 (fun (m0 : nat) (IHle : m ≤ m0) (IHle : n ≤ m0 + 1) =>	14 (fun (m0 : nat) (IHle : m ≤ m0) (IHle : n ≤ m0) =>
15 le_S n (m0 + 1) IHle)	15 le_S n m0 IHle)
16 p <span style="float: right;">(* p *)</span>	16 p <span style="float: right;">(* p *)</span>
17 H0	17 H0

**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.

for a reusable patch. The other components modify the candidate to try to produce a patch. To motivate this workflow, consider using PUMPKIN to search the proofs in Figure 2 for a patch between conclusions. 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):

```
 $\forall n\ m\ p, n \leq m \rightarrow m \leq p \rightarrow n \leq p \rightarrow n \leq p + 1$ 
```

It then breaks each inductive proof into cases and determines an intermediate goal type for the candidate. In the base case, for example, 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)
  :  $\forall n\ m\ p, n \leq m \rightarrow m \leq p \rightarrow$ 
     $n \leq m \rightarrow n \leq m + 1$ 
```

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`). 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)
  :  $\forall n0\ n\ m\ p, n \leq m \rightarrow m \leq p \rightarrow$ 
     $n \leq n0 \rightarrow n \leq n0 + 1$ 
```

PUMPKIN then *specializes* this candidate to `p` (line 16), the argument to the conclusion of `le_ind`. This produces the patch:

```
patch n m p H0 H1 :=
  (fun (H : n ≤ p) => le_plus_trans n p 1 H)
  :  $\forall n\ m\ p, n \leq m \rightarrow m \leq p \rightarrow$ 
     $n \leq p \rightarrow n \leq p + 1$ 
```

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 simple example uses only three components. The other components help turn candidates into patches in similar ways. We discuss all five components in more detail in the rest of this section.

**Semantic 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):

```
le_plus_trans n m 1 H : n ≤ m + 1
H : n ≤ m
```

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:

```
fun (H : n ≤ m) => 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>2</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:

```
... (IHle : n ≤ m0 + 1) ...
... (IHle : n ≤ m0) ...
```

However, the type of `IHle` changes for *any* two inductive proofs over `le` with different conclusions. A syntactic

<sup>2</sup>Since this is a simple example, replaying an existing tactic happens to work. There are additional examples in the repository (Cex.v).

differencing component may identify this change as a candidate. Our semantic differencing component knows that it can ignore this change.

**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 it to a hint database [3], 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. ✓
```

to this proof:

```
intros. rewrite rev_app_distr. apply new. ✓
```

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 sub-problems. It is also necessary to invert certain terms. Consider inverting an arbitrary sequence of two rewrites:

```
t := eq_ind_r G ... (eq_ind_r F ...)
```

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

```
g := eq_ind_r G ...
f := eq_ind_r F ...
t := g ∘ f
```

The inverse of  $t$  is the following:

```
t-1 := f-1 ∘ 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.

### 3.2 Composing Components

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: if direction is forwards then
7:   find_patch(term', term, backwards) for candidates
8:   factor and invert candidates
9:   if there are patches then return patches
10: 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:

```
∀ x, H x -> P x
∀ x, H x -> P' x
```

PUMPKIN searches for a patch with this type:

```
∀ x, H x -> P' x -> 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
```

---

Section 4 describes real-world examples that demonstrate more configurations.

```
Record int : Type :=
  mkint { intval: Z; intrange: 0 <= intval < modulus }.
```

```
Record int : Type :=
  mkint { intval: Z; intrange: -1 < intval < modulus }.
```

Figure 3. Old (left) and new (right) definitions of `int` in CompCert.

## 4 Case Studies: Changes in the Wild

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

1. **Updating definitions** within a project (CompCert, Section 4.1)
2. **Porting definitions** between libraries (Software Foundations, Section 4.2)
3. **Updating proof assistant versions** (Coq Standard Library, Section 4.3)

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 4.1 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 preserved shape, like modifying constructors. We discuss supporting these features and whether they may necessitate other components in Section 7.

### 4.1 Updating Definitions

Coq programmers sometimes make changes to definitions that break proofs within the same project. To emulate this use case, we identified a CompCert commit [44] with a breaking change to `int` (Figure 3). 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.✓
```

#### 4.1.1 Core Components

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>3</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
- 

<sup>3</sup>If PUMPKIN finds just one implication, it returns that.



```

Fixpoint bin_to_nat (b : bin) : nat := match b with
| B0 => 0
| B2 b' => 2 * (bin_to_nat b')
| B21 b' => 1 + 2 * (bin_to_nat b')
end.

```

```

Fixpoint bin_to_nat (b : bin) : nat := match b with
| B0 => 0
| B2 b' => (bin_to_nat b') + (bin_to_nat b')
| B21 b' => S ((bin_to_nat b') + (bin_to_nat b'))
end.

```

Figure 4. Definitions of `bin_to_nat` for Users A (left) and B (right).

## 4.2 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 [10, 14] to an exercise in Software Foundations to each use the other solution’s definition of the fixpoint `bin_to_nat` (Figure 4). 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. ✗

```

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>4</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

<sup>4</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.

changes do not provide enough information to determine this; the tool must have a semantic understanding of the terms.

### 4.2.1 Core Components

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:

```

∀ ..., P (f ...)
∀ ..., P (f' ...)

```

for an isomorphism that corresponds to the change:

```

∀ P, P x -> P x'
∀ P, P x' -> P x

```

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

## 4.3 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 [45]. The commit changed the definition of `divide` prior to the release of Coq 8.4 (Figure 5). The change broke 46 proofs in the standard library. We used PUMPKIN to find an isomorphism that corresponds

**Definition** `divide p q := ∃ r, p * r = q.`

**Definition** `divide p q := ∃ r, q = r * p.`

**Figure 5.** Old (left) and new (right) definitions of `divide` in Coq.

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:`  
 $\forall a\ b, b \neq 0 \rightarrow (a \bmod b = 0 \leftrightarrow (\text{divide } b\ a)).$

**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:

$\forall r\ p\ q, p * r = q \rightarrow q = r * p$   
 $\forall r\ p\ q, q = r * p \rightarrow p * r = q$

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

**Theorem** `Zmod_divides:`  
 $\forall a\ b, b > 0 \rightarrow (a \bmod b = 0 \leftrightarrow \exists c, a = b * c).$   
**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. ✓`

#### 4.3.1 Core Components

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:

`... (C (P x)) ...`  
`... (C (P' x)) ...`

for an isomorphism between the arguments:

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

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 any need for induction; a patch between constructors does not. For example, for this case study, 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.

## 5 Inside the Core

We have shown that the core components are useful for finding proof patches. This section describes our implementation of the core components (Section 5.1) and the procedure (Section 5.2) so that it is possible to implement them in other systems. While our system is a very early prototype under active development, we have made the current source code available on Github.<sup>5</sup> Our prototype has no significant impact on the trusted computing base (Section 5.3).

### 5.1 Inside the Components

The interested reader can follow along with the descriptions of the components in the repository.

#### 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>6</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.

<sup>5</sup><http://github.com/uwplse/PUMPKIN-PATCH/tree/cpp18>

<sup>6</sup>These trees are inspired by categorical models of dependent type theory [36].

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 6) 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 only 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.

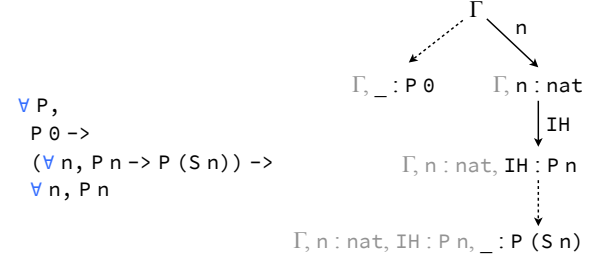
The tool always maintains pointers to easily switch between the tree and AST representations of the terms. This representation enables extensibility; we discuss a broad range of extensions in Sections 6.2 and 7.

### 5.1.2 Patch Specialization

Specialization (`specialize.ml`) takes a patch candidate and some arguments, all of which are terms in Coq's internal representation. It applies the candidate to the arguments, then it  $\beta$ -reduces [25] 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.

### 5.1.3 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



**Figure 6.** 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.

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* occurrences of terms convertible to an argument or function (we show an example in Section 6.2), or the argument or function may not show up in the term 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. It moves on to more complicated strategies only if the first strategy fails. This design makes abstraction simple to extend with new



strategies and simple to call with different strategies for different classes of changes; we identify new strategies in Section 6.2.

#### 5.1.4 Patch Inversion

Patch inversion (`inverting.ml`) exploits symmetry to try to reverse the conclusions of a patch. It factors the patch using the factoring component, then calls the primitive inversion function on each factor. It then 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. In the future, it may be possible to use this swapping functionality in a reflexive application of the induction principle to infer symmetry properties like `eq_sym` for other types. The component can then use these properties to generate reverse induction principles like `eq_ind_r`. We further detail inversion in Section 6.2.

#### 5.1.5 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 -> y` for some `y`. It then recursively searches for a term with type `y -> 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.

## 5.2 Inside the Procedure

The implementation (`patcher.ml4`) of the procedure from Section 3.2 starts with a preprocessing step which compiles the proof terms to trees (like the tree in Figure 6). 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`. 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 4.1, 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 3), 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.

## 5.3 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 [30]. 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 [40]. Furthermore, Coq type-checks terms that plugins produce. Since PUMPKIN does not modify the type checker, it cannot produce an ill-typed term.

## 6 Testing Boundaries

In the case studies in Section 4, we showed how the core components are useful for real scenarios. In this section, we explore the boundary between what the PUMPKIN prototype can and cannot handle. It is precisely this boundary that informs us how to improve the implementations of the core components.

To evaluate this boundary, we tested the core components of PUMPKIN on a suite of 50 pairs of proofs (Section 6.1). We designed 11 of these pairs to succeed, then modified their proofs to produce the remaining 39 pairs that try to stress the core functionality of the tool. We learned the following from the pairs that tested PUMPKIN's limitations:

### 1. The failed pairs drive improvements.

PUMPKIN failed on 17 of 50 pairs. These pairs tell us how to improve the core. (Section 6.2)

### 2. The pairs unearth abstraction strategies.

PUMPKIN produced an exponential number of candidates in 5 of 50 pairs. New abstraction strategies can dramatically reduce the number of candidates. (Section 6.2)

### 3. Pumpkin was fast, and it can be faster.

The slowest successful patch took 48 ms. The slowest failure took 7 ms. Simple changes can make PUMPKIN more efficient. (Section 6.3)

```

fun n m p (H : n <= m) (H0 : m <= p) =>
  le_S n p (* ... proof of stronger lemma *)
: ∀ n m p, n <= m -> m <= p -> n <= S p

fun n m p (H : n <= m) (H0 : m <= p) =>
  le_plus_trans n p 1 (* ... proof of stronger lemma *)
: ∀ n m p, n <= m -> m <= p -> n <= p + 1

```

Figure 7. Two proof terms `old` (left) and `new` (right) that contain the same proof of a stronger lemma.

### 6.1 Patch Generation Suite

We wrote a suite<sup>7</sup> of 50 pairs of proofs. We wrote these proofs ourselves since searching for proof patches is a new domain, so there was no existing benchmark suite to work with. We used the following methodology:

1. Choose theorems `old` and `new`
2. Write similar inductive proofs of `old` and `new`
3. Modify the proof of `old` to produce more pairs
4. Search for patches from `new` to `old`
5. If possible, search for patches from `old` to `new`

In total, we chose 11 pairs of theorems `old` and `new`, and we wrote 50 pairs of proofs of those theorems.

For example, one pair of theorems `old` and `new` was a simplification of the auxiliary lemmas that we encountered in the case study in Section 4.2. For the first proof of `old`, we added a rewrite, like in the case study:

```
rewrite <- plus_n0. rewrite -> plus_comm.
```

For the second proof of `old`, we commuted the rewrites:

```
rewrite -> plus_comm. rewrite <- plus_n0.
```

We then searched for patches in both directions, since the conclusions of `old` and `new` were propositionally equal.

Our goal was to determine what changes to proofs stress the core components and how to use that information to drive improvements. We focused on differences in conclusions, the most supported configuration in our tool. Since PUMPKIN operates at the term level, we removed redundant proof terms, even if they were produced by different tactics. We controlled the first pair of proofs of each pair of theorems for features we had not yet implemented, like nested induction, changes in hypotheses, and abstracting `omega` terms. These features sometimes showed up in later proofs (for example, after moving a rewrite); we kept these proofs in the patch generation suite, since isolated changes to supported proofs that introduce unsupported features can inform future improvements.

### 6.2 Three Challenges

PUMPKIN found patches for 33 of the 50 pairs. 28 of the 33 successes did not stress PUMPKIN at all: PUMPKIN found the correct candidate immediately and was able to abstract it in one try. The pairs that PUMPKIN failed to patch and the successful pairs that stressed abstraction reveal key information about how to improve the core components. We walk through three examples. Future

tools can use these as challenge problems to improve upon PUMPKIN.

**A Challenge for Differencing** For one pair of proofs of theorems with propositionally equal conclusions (Figure 7), the differencing component failed to find candidates in either direction. These proofs both contain the same proof of a stronger lemma; PUMPKIN found patches from this lemma to both `old` and `new`, but it was unable to find a patch between `old` and `new`. A patch may show up deep in the difference between `le_plus_trans` and `le_S`, but even if we  $\delta$ -reduce (unfold the definition of [25]) `le_plus_trans`, this is not obvious:

```

le_plus_trans n m p (H : n <= m) :=
  (fun lemma : m <= m + p =>
    trans_contra_inv_impl_morphism
      PreOrder_Transitive
        (m + p)
        m
        lemma)
    (le_add_r m p)
  H

```

This points to two difficulties in finding patches: Knowing when to  $\delta$ -reduce terms is difficult; exploring the appropriate time for reduction may produce patches for pairs that PUMPKIN currently cannot patch. Furthermore, finding patches is more challenging when neither theorem has a conclusion that is as strong as possible.

**A Challenge for Inversion** For one pair of proofs with propositionally equal conclusions, PUMPKIN found a patch in one direction, but failed to invert it:

```

fun n m p (_ : n <= m) (_ : m <= p) (H1 : n <= p) =>
  gt_le_S n (S p) (le_lt_n_Sm n p H1)
: ∀ n m p, n <= m -> m <= p -> n <= p -> S n <= S p

```

The inversion component was unable to invert this term, even though an inverse does exist. To invert this, the component needs to know to  $\delta$ -reduce `gt_le_S`:

```

gt_le_S n m :=
  (fun (H : ∀ n0 m0, n0 < m0 -> S n0 <= m0) =>
    H n m)
  ...
: ∀ n m, n < m -> S n <= m

It then needs to swap the hypothesis with the conclusion in H to produce the inverse:

gt_le_S-1 n m :=
  (fun (H : ∀ n0 m0, S n0 <= m0 -> n0 < m0) =>
    H n m)
  ...
: ∀ n m, S n <= m -> n < m

```

<sup>7</sup><http://github.com/uwplse/PUMPKIN-PATCH/blob/cpp18/plugin/coq/Variants.v>

Inversion currently swaps subterms when it is not aware of any symmetry properties about the inductive type. However, it does not know when to  $\delta$ -reduce functions. Furthermore, there are many possible subterms to swap; for inversion to know to only swap the subterms of  $n$ , it must have a better understanding of the structure of the term. Both of these are ways to improve inversion.

**A Challenge for Abstraction** Abstraction produced an exponential number of candidates when abstracting a patch candidate with this type:

```

∀ n n0,
  (fun m => n <= max m n0) n ->
  (fun m => n <= max n0 m) n

```

The goal was to abstract by  $n$  and produce a patch with this type:

```

∀ m0 n n0, n <= max m0 n0 -> n <= max n0 m0

```

The difficulty was in determining which occurrences of  $n$  to abstract. The component needed to abstract only the highlighted occurrences:

```

fun n n0 (H0 : n <= max n0 n) =>
  @eq_ind_r nat
    (max n0 n)
    (fun n1 => n <= n1)
    H0
    (max n n0)
    (max_comm n n0)

```

The simplest abstraction strategy failed, and a more complex strategy succeeded only after producing exponentially many candidates. While this did not have a significant impact on time, this case gives rise to a new class of abstraction strategies: semantics-aware abstraction. In this case, we know from the type of the candidate and the type of `eq_ind_r` that these two hypothesis types are equivalent (similarly for the conclusion types):

```

(fun m => n <= max m n0) n
(fun n1 => n <= n1) (max n0 n)

```

The tool can search recursively for patches to find two patches that bridge the two equivalent types:

```

p1 := fun n => max n0 n
p2 := fun n => max n n0

```

Then the candidate type is exactly this:

```

∀ n n0,
  (fun n1 => n <= n1) (p2 n) ->
  (fun n1 => n <= n1) (p1 n)

```

Abstraction should thus abstract the highlighted subterms and the terms that have types constrained by those subterms. This produces a patch in one candidate:

```

fun m0 n n0 (H0 : n <= max n0 m0) =>
  @eq_ind_r nat
    (max n0 m0) (* p1 m0 *)
    (fun n1 => n <= n1) (* P *)
    H0 (* : P (p1 m0) *)
    (max m0 n0) (* p2 m0 *)
    (max_comm m0 n0) (* : p1 m0 = p2 m0 *)

```

This same strategy would find a patch for one of the pairs that PUMPKIN failed to abstract. This is a natural future direction for abstraction.

### 6.3 Performance

PUMPKIN performed well for all pairs. The slowest success took 48 ms.<sup>8</sup> When PUMPKIN failed, it failed fast. The slowest failure took 7 ms. While we find this promising, proof terms were small ( $\leq 67$  LOC); we leave evaluating performance on larger terms to future work.

PUMPKIN was slowest when the patch showed up inverted in the difference of proofs, since PUMPKIN had to search twice, once in each direction. A future procedure may determine which direction to search first ahead of time; proof term size may be a simple heuristic for this.

## 7 Conclusions & Future Work

Our vision is a future of proof automation in ITP that adapts proofs to breaking changes. This will unify the spirit of ITP—collaboration between the programmer and the tool—with the realities of modern proof engineering: Verification projects are large, specifications evolve over time, and dependencies change and break backward-compatibility. Too much of the burden of change rests on the programmer; not enough rests on the tool.

We conclude with a discussion of improvements that can help bring this vision to life. These improvements are driven by our experiences using the PUMPKIN prototype and by conversations within the ITP community.

**Supporting structural changes.** Coq programmers often make structural changes. It is common, for example, to add new hypotheses, constructors, or parameters to a type. The ideal tool should find patches for these changes. Existing work in proof reuse [17, 54] and type-directed diffing [52] may help guide these improvements.

**Exploring new components.** The core components of PUMPKIN are critical to searching for patches in Coq, but they may not be sufficient for the ideal tool. While we find the flexibility of these components promising thus far, in implementing new features, we may discover new components. For example, patches for certain structural changes may be program transformations as opposed to terms; it remains to be seen whether supporting such patches will reveal new components.

**Modeling diverse proof styles.** Coq programmers use diverse proof styles; the ideal tool should support many different styles. Proofs in decidable domains that apply `dec_not_not` pose difficulties for abstraction and inversion; the ideal tool should support these. PUMPKIN has limited support for changes in hypotheses, fixpoints,

<sup>8</sup>i7-4790K, at 4.00 GHz, 32 GB RAM

constructors, pattern matching, and nested induction; the ideal tool should implement these features.

**Improving user experience.** The ideal proof patching tool should be natural for programmers to use. A future patching tool can produce tactics from the patches it finds, that way programmers can remove references to old specifications. A future patching tool can integrate with an IDE (such as Proof General [8]) or continuous integration (CI) system (such as Travis [9]) to suggest patches at natural steps in the development process.

**Handling version updates.** Updating Coq versions is an ideal use case for finding proof patches: When the client updates Coq, a tool can automatically search commits to the standard library or to `coq-contribs` for patches. In this case, the example comes from the Coq developers, not from the library client—the client never has to look at the changes that the Coq developers make. The ideal tool should fully support updates in Coq versions. PUMPKIN can patch certain changes in the standard library, but it does not yet search for those changes automatically and determine which configuration to use depending on what has changed. Furthermore, as proof assistant versions change, so may the AST and the plugin interface. To fully support version updates, the tool should support different language versions.

**Isolating changes.** Programmers sometimes make multiple changes to a verification project in the same commit; the ideal tool should break down large changes into small, isolated changes to find patches. This will help in finding benchmarks, supporting library and version updates, and integrating with CI systems. We may draw on work in change and dependency management [11, 22, 38] to identify changes, then use the factoring component to break these changes into smaller parts.

**Supporting other proof assistants.** Coq is just one of many proof assistants; ideal tools should support different proof assistants. While PUMPKIN focuses on Coq, the underlying concepts extend to other proof assistants with constructive logics (for example, Lean [6] and Agda [1]). Proof assistants with non-constructive logics (for example, Isabelle/HOL [5]) may benefit from a different approach; this is similar to the problem of finding patches for proofs in decidable domains in Coq, since classical properties provably hold for decidable propositions [7].

**Applying patches.** The ideal tool should not only find patches, but also apply the patches it finds automatically to fix broken proofs. In some cases, this may be as simple as adding hints to a hint database, so that proofs go through with no changes. However, hint databases in Coq cannot support certain terms [3], and adding too

many hints may impede performance. More generally, we can integrate PUMPKIN with a transfer tactic [37, 65], which is a perfect fit: Transfer tactics automatically adapt proofs between isomorphic types and implications, but they do not find these functions; PUMPKIN finds these functions, but it does not apply them.

## 8 Related Work

Our work builds upon prior research in proof reuse, proof automation, proof engineering, refactoring, differencing & incremental computation, programming by example, and program repair.

**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 [15] or through an automatic method [47]. This is later generalized and implemented in Isabelle [37] and Coq [61, 65]; 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 [17] adapts proof obligations to structural changes in inductive types. Later work [54] 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 [31, 39, 58] 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.

**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` [24] can make proofs more resilient to changes. Hammers like Isabelle’s `sledgehammer` [56] can make proofs agnostic to some low-level changes. Recent work [26] 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 [23, 24] can help prevent low-level maintenance tasks. Writing these tactics requires good engineering [33] 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 [20] 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 [51] may help improve these components. Lean [60] 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 [19] repairs faulty physics theories in the context of a classical higher-order logic (HOL); there is preliminary work extending this style of repair to mathematical proofs. Knowledge-sharing methods [32] can adapt some proofs across different representations of HOL. These complementary approaches may guide extensions to support decidable domains and classical logics.

**Proof Engineering** Existing proof engineering work addresses brittleness by planning for changes [64] and designing theorems and proofs that make maintenance less of an issue. Design principles for specific domains (such as formal metatheory [13, 28, 29]) can make verification more tractable. CertiKOS [34] introduces the idea of a deep specification to ease verification of large systems. Ornaments [27, 63] separate the computational and logical components of a datatype, and may make proofs more resilient to datatype changes. 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.

There is a small body of work on change and dependency management for verification, both to evaluate impact of potential changes and maximize reuse [11, 38] and to optimize build performance [22]. These approaches may help in isolating changes, which is necessary to identify future benchmarks, integrate with CI systems, and fully support version updates.

**Refactoring** Our approach is close in spirit to refactoring [50]. The Haskell refactoring tool HaRe [4] automatically lifts definitions, and may be useful for improving abstraction. There is a growing body of work on refactoring in the context of ITP [18, 62]. The IDE CoqPIE [59] and the verification platform Why3 [16] can both adapt Coq proofs to simple syntactic changes.

It may be possible to use our lemma factoring component to improve proof refactoring tools. Proof refactoring tools are semantics-preserving; unlike these tools, our approach handles semantic changes.

**Differencing & Incremental Computation** Existing work in differencing and incremental computation may help improve our semantic differencing component. Type-directed diffing [52] finds differences in algebraic data types. Semantics-based change impact analysis [12] models semantic differences between documents. Differential assertion checking [41] analyzes different versions of a program for relative correctness with respect to a specification. Incremental  $\lambda$ -calculus [21] introduces a general model for program changes. All of these may be useful for improving semantic differencing.

**Programming by Example** Our approach generalizes an example that the programmer provides. This is similar to programming by example, a subfield of program synthesis [35]. This field addresses different challenges in different logics, but may drive solutions to similar problems in a dependently typed language.

**Program Repair** Adapting proofs to changes is essentially program repair for dependently typed languages. Program repair tools for languages with non-dependent type systems [42, 46, 48, 53, 57] 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.

## 9 Acknowledgements

We thank Valentin Robert and Zach Tatlock for valuable conversations about proof patching. We thank Leonardo de Moura and the UW PLSE lab for helpful suggestions. We thank the students and lecturers from the first DeepSpec Summer School for motivating future work. We thank Dan Licata for pointers in understanding the mathematical properties of inductive types. This material is based upon work supported by the National Science Foundation Graduate Research Fellowship under Grant No. DGE-1256082. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

## References

- [1] 2017. Agda. (2017). <http://wiki.portal.chalmers.se/agda/pmwiki.php>
- [2] 2017. Coq. (2017). <http://coq.inria.fr/>
- [3] 2017. Coq Reference Manual, Section 8.9: Controlling Automation. (2017). <http://coq.inria.fr/refman/tactics.html>
- [4] 2017. HaRe: The Haskell Refactoring Tool. (2017). <http://github.com/RefactoringTools/HaRe>



- [5] 2017. Isabelle/HOL: A Proof Assistant for Higher-Order Logic. (2017). <http://isabelle.in.tum.de/doc/tutorial.pdf>
- [6] 2017. Lean Theorem Prover. (2017). <http://github.com/leanprover/lean>
- [7] 2017. Library Coq.Logic.Decidable. (2017). <http://coq.inria.fr/library/Coq.Logic.Decidable.html>
- [8] 2017. Proof General. (2017). <http://proofgeneral.github.io/>
- [9] 2017. Travis CI. (2017). <http://travis-ci.org/>
- [10] User A. 2017. Software Foundations Solution. (2017). <http://github.com/blindFS/Software-Foundations-Solutions>
- [11] Serge Autexier, Dieter Hutter, and Till Mossakowski. 2010. Verification, Induction Termination Analysis. Springer-Verlag, Berlin, Heidelberg, Chapter Change Management for Heterogeneous Development Graphs, 54–80. <http://dl.acm.org/citation.cfm?id=1986659.1986663>
- [12] Serge Autexier and Normen Müller. 2010. Semantics-based Change Impact Analysis for Heterogeneous Collections of Documents. In *Proceedings of the 10th ACM Symposium on Document Engineering (DocEng '10)*. ACM, New York, NY, USA, 97–106. DOI:<http://dx.doi.org/10.1145/1860559.1860580>
- [13] Brian Aydemir, Arthur Charguéraud, Benjamin C. Pierce, Randy Pollack, and Stephanie Weirich. 2008. Engineering Formal Metatheory. In *Proceedings of the 35th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '08)*. ACM, New York, NY, USA, 3–15. DOI:<http://dx.doi.org/10.1145/1328438.1328443>
- [14] User B. 2017. Software Foundations Solution. (2017). [http://github.com/marshall-lee/software\\_foundations](http://github.com/marshall-lee/software_foundations)
- [15] Gilles Barthe and Olivier Pons. 2001. 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)*. Springer-Verlag, London, UK, 57–71. <http://dl.acm.org/citation.cfm?id=646793.704711>
- [16] François Bobot, Jean-Christophe Filliâtre, Claude Marché, Guillaume Melquiond, and Andrei Paskevich. 2013. Preserving User Proofs Across Specification Changes. In *Fifth Working Conference on Verified Software: Theories, Tools and Experiments*, Ernie Cohen and Andrey Rybalchenko (Eds.), Vol. 8164. Springer, Atherton, United States, 191–201. <https://hal.inria.fr/hal-00875395>
- [17] Olivier Boite. 2004. Proof Reuse with Extended Inductive Types. In *Theorem Proving in Higher Order Logics: 17th International Conference, TPHOLs 2004, Park City, Utah, USA, September 14-17, 2004. Proceedings*, Konrad Slind, Annette Bunker, and Ganesh Gopalakrishnan (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 50–65. DOI:[http://dx.doi.org/10.1007/978-3-540-30142-4\\_4](http://dx.doi.org/10.1007/978-3-540-30142-4_4)
- [18] Timothy Bourke, Matthias Daum, Gerwin Klein, and Rafal Kolanski. 2012. Challenges and Experiences in Managing Large-Scale Proofs. In *Conferences on Intelligent Computer Mathematics (CICM) / Mathematical Knowledge Management*, Makarius Wenzel (Ed.). Springer, Bremen, Germany, 32–48.
- [19] Alan Bundy. 2013. The interaction of representation and reasoning. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 469, 2157 (2013). DOI:<http://dx.doi.org/10.1098/rspa.2013.0194>
- [20] Alan Bundy, David Basin, Dieter Hutter, and Andrew Ireland. 2005. *Rippling: Meta-Level Guidance for Mathematical Reasoning*. Cambridge University Press, New York, NY, USA.
- [21] Yufei Cai, Paolo G. Giarrusso, Tillmann Rendel, and Klaus Ostermann. 2014. 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)*. ACM, New York, NY, USA, 145–155. DOI:<http://dx.doi.org/10.1145/2594291.2594304>
- [22] Ahmet Celik, Karl Palmkog, and Milos Gligoric. 2017. iCoq: Regression Proof Selection for Large-scale Verification Projects. In *Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering (ASE 2017)*. IEEE Press, Piscataway, NJ, USA, 171–182. <http://dl.acm.org/citation.cfm?id=3155562.3155588>
- [23] Adam Chlipala. 2013. 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)*. ACM, New York, NY, USA, 391–402. DOI:<http://dx.doi.org/10.1145/2500365.2500592>
- [24] Adam Chlipala. 2013. *Certified Programming with Dependent Types - A Pragmatic Introduction to the Coq Proof Assistant*. MIT Press. <http://mitpress.mit.edu/books/certified-programming-dependent-types>
- [25] Adam Chlipala. 2017. Library Equality. (2017). <http://adam.chlipala.net/cpdt/html/Equality.html>
- [26] Lukasz Czapka and Cezary Kaliszyk. 2017. Hammer for Coq: Automation for Dependent Type Theory. (2017). <http://cl-informatik.uibk.ac.at/cek/coqhammer/>
- [27] Pierre-Évariste Dagand. 2017. The essence of ornaments. *J. Funct. Program.* 27 (2017), e9. DOI:<http://dx.doi.org/10.1017/S0956796816000356>
- [28] Benjamin Delaware, Bruno C. d. S. Oliveira, and Tom Schrijvers. 2013. Meta-theory à La Carte. In *Proceedings of the 40th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '13)*. ACM, New York, NY, USA, 207–218. DOI:<http://dx.doi.org/10.1145/2429069.2429094>
- [29] Benjamin Delaware, Steven Keuchel, Tom Schrijvers, and Bruno C.d.S. Oliveira. 2013. Modular Monadic Meta-theory. In *Proceedings of the 18th ACM SIGPLAN International Conference on Functional Programming (ICFP '13)*. ACM, New York, NY, USA, 319–330. DOI:<http://dx.doi.org/10.1145/2500365.2500587>
- [30] Maxime Dénès. 2017. Coq 8.7 beta 1 is out. (2017). <http://coq.inria.fr/news/137.html>
- [31] Amy Felty and Douglas Howe. 1994. Generalization and reuse of tactic proofs. In *Logic Programming and Automated Reasoning: 5th International Conference (LPAR '94)*, Frank Pfenning (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 1–15. DOI:[http://dx.doi.org/10.1007/3-540-58216-9\\_25](http://dx.doi.org/10.1007/3-540-58216-9_25)
- [32] Thibault Gauthier and Cezary Kaliszyk. 2014. Matching concepts across HOL libraries. In *CICM '14 (LNCS)*, Stephen Watt, James Davenport, Alan Sexton, Petr Sojka, and Josef Urban (Eds.), Vol. 8543. Springer Verlag, 267–281. DOI:[http://dx.doi.org/10.1007/978-3-319-08434-3\\_20](http://dx.doi.org/10.1007/978-3-319-08434-3_20)
- [33] Georges Gonthier, Beta Ziliani, Aleksandar Nanevski, and Derek Dreyer. 2011. How to Make Ad Hoc Proof Automation Less Ad Hoc. In *Proceedings of the 16th ACM SIGPLAN International Conference on Functional Programming (ICFP '11)*. ACM, New York, NY, USA, 163–175. DOI:<http://dx.doi.org/10.1145/2034773.2034798>
- [34] Ronghui Gu, Zhong Shao, Hao Chen, Xiongnan (Newman) Wu, Jieung Kim, Vilhelm Sjöberg, and David Costanzo. 2016. CertiKOS: An Extensible Architecture for Building Certified Concurrent OS Kernels. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*.

- USENIX Association, GA, 653–669. <https://www.usenix.org/conference/osdi16/technical-sessions/presentation/gu>
- [35] Sumit Gulwani, Oleksandr Polozov, and Rishabh Singh. 2017. Program Synthesis. *Foundations and Trends in Programming Languages* 4, 1-2 (2017), 1–119. DOI:<http://dx.doi.org/10.1561/25000000010>
- [36] Martin Hofmann. 1997. Syntax and Semantics of Dependent Types. In *Semantics and Logics of Computation*. Cambridge University Press, 79–130.
- [37] Brian Huffman and Ondřej Kunčar. 2013. Lifting and Transfer: A Modular Design for Quotients in Isabelle/HOL. In *Certified Programs and Proofs: Third International Conference (CPP 2013)*, Georges Gonthier and Michael Norrish (Eds.). Springer International Publishing, Cham, 131–146. DOI:[http://dx.doi.org/10.1007/978-3-319-03545-1\\_9](http://dx.doi.org/10.1007/978-3-319-03545-1_9)
- [38] D. Hutter. 2000. Management of change in structured verification. In *ASE 2000*. 23–31. DOI:<http://dx.doi.org/10.1109/ASE.2000.873647>
- [39] Einar Broch Johnsen and Christoph Lüth. 2004. Theorem Reuse by Proof Term Transformation. In *Theorem Proving in Higher Order Logics: 17th International Conference (TPHOLs 2004)*, Konrad Slind, Annette Bunker, and Ganesh Gopalakrishnan (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 152–167. DOI:[http://dx.doi.org/10.1007/978-3-540-30142-4\\_12](http://dx.doi.org/10.1007/978-3-540-30142-4_12)
- [40] Matej Kosik. 2017. Coq Pull Request # 652: Put all plugins behind an “API”. (2017). <http://github.com/coq/coq/pull/652>
- [41] Shuvendu Lahiri, Kenneth McMillan, , and Chris Hawblitzel. 2013. Differential Assertion Checking, In *Foundations of Software Engineering (FSE’13)*. (August 2013). <https://www.microsoft.com/en-us/research/publication/differential-assertion-checking-2/>
- [42] Xuan-Bach D. Le, Duc-Hiep Chu, David Lo, Claire Le Goues, and Willem Visser. 2017. 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)*. ACM, New York, NY, USA, 593–604. DOI:<http://dx.doi.org/10.1145/3106237.3106309>
- [43] Xavier Leroy. 2006. Formal certification of a compiler backend, or: programming a compiler with a proof assistant. In *33rd ACM symposium on Principles of Programming Languages*. ACM Press, 42–54. <http://gallium.inria.fr/~xleroy/publi/compiler-certif.pdf>
- [44] Xavier Leroy. 2013. Commit to CompCert: lib/Integers.v. (2013). <http://github.com/AbsInt/CompCert/commit/6f3225b0623b9c97eed7d40ddc320b08c79c6518>
- [45] letouzey. 2011. Commit to coq: change definition of divide (compat with Znumtheory). (2011). <http://github.com/coq/coq/commit/81c4c8bc418cdf42cc88249952dbba465068202c>
- [46] Fan Long and Martin Rinard. 2016. Automatic Patch Generation by Learning Correct Code. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL ’16)*. ACM, New York, NY, USA, 298–312. DOI:<http://dx.doi.org/10.1145/2837614.2837617>
- [47] Nicolas Magaud and Yves Bertot. 2002. Changing Data Structures in Type Theory: A Study of Natural Numbers. In *Types for Proofs and Programs: International Workshop (TYPES 2000)*, Paul Callaghan, Zhaoxue Luo, James McKeenna, and Robert Pollack (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 181–196. DOI:[http://dx.doi.org/10.1007/3-540-45842-5\\_12](http://dx.doi.org/10.1007/3-540-45842-5_12)
- [48] Sergey Mechtaev, Jooyong Yi, and Abhik Roychoudhury. 2016. Angelix: Scalable Multiline Program Patch Synthesis via Symbolic Analysis. In *Proceedings of the 38th International Conference on Software Engineering (ICSE ’16)*. ACM, New York, NY, USA, 691–701. DOI:<http://dx.doi.org/10.1145/2884781.2884807>
- [49] Guillaume Melquiond. 2017. Commit to coq: Make IZR use a compact representation of integers. (2017). <http://github.com/coq/coq/commit/a4a76c253474ac4ce523b70d0150ea5dcf546385>
- [50] Tom Mens and Tom Tourwé. 2004. A Survey of Software Refactoring. *IEEE Trans. Softw. Eng.* 30, 2 (Feb. 2004), 126–139. DOI:<http://dx.doi.org/10.1109/TSE.2004.1265817>
- [51] Dale Miller and Gopalan Nadathur. 2012. *Programming with Higher-Order Logic* (1st ed.). Cambridge University Press, New York, NY, USA.
- [52] Victor Cacciari Miraldo, Pierre-Évariste Dagand, and Wouter Swierstra. 2017. Type-directed Diffing of Structured Data. In *Proceedings of the 2Nd ACM SIGPLAN International Workshop on Type-Driven Development (TyDe 2017)*. ACM, New York, NY, USA, 2–15. DOI:<http://dx.doi.org/10.1145/3122975.3122976>
- [53] Martin Monperrus. 2017. Automatic Software Repair: a Bibliography. *ACM Computing Surveys* (2017). <https://hal.archives-ouvertes.fr/hal-01206501/file/survey-automatic-repair.pdf>
- [54] Anne Mulhern. 2006. Proof Weaving. In *In Proceedings of the First Informal ACM SIGPLAN Workshop on Mechanizing Metatheory*.
- [55] Karl Palmskog. 2017. Commit to verdi-raft: Port to Coq 8.6. (2017). <http://github.com/uwplse/verdi-raft/pull/43/files>
- [56] L. C. Paulson and J. C. Blanchette. 2012. Three years of experience with Sledgehammer, a practical link between automatic and interactive theorem provers. In *International Workshop on the Implementation of Logics (IWIL 2010) (EPiC Series)*, G. Sutcliffe, S. Schulz, and E. Ternovska (Eds.), Vol. 2. EasyChair, 1–11.
- [57] Yu Pei, Carlo A. Furia, Martin Nordio, and Bertrand Meyer. 2014. Automatic Program Repair by Fixing Contracts. In *Proceedings of the 17th International Conference on Fundamental Approaches to Software Engineering - Volume 8411*. Springer-Verlag New York, Inc., New York, NY, USA, 246–260. DOI:[http://dx.doi.org/10.1007/978-3-642-54804-8\\_17](http://dx.doi.org/10.1007/978-3-642-54804-8_17)
- [58] Olivier Pons. 2000. Generalization in Type Theory Based Proof Assistants (*TYPES ’00*). 217–232.
- [59] Kenneth Roe and Scott Smith. 2016. CoqPIE: An IDE Aimed at Improving Proof Development Productivity. In *Interactive Theorem Proving: 7th International Conference, ITP 2016, Nancy, France, August 22-25, 2016, Proceedings*, Jasmin Christian Blanchette and Stephan Merz (Eds.). Springer International Publishing, Cham, 491–499. DOI:[http://dx.doi.org/10.1007/978-3-319-43144-4\\_32](http://dx.doi.org/10.1007/978-3-319-43144-4_32)
- [60] Daniel Selsam and Leonardo de Moura. 2016. Congruence Closure in Intensional Type Theory. In *Automated Reasoning: 8th International Joint Conference (IJCAR 2016)*, Nicola Olivetti and Ashish Tiwari (Eds.). Springer International Publishing, Cham, 99–115. DOI:[http://dx.doi.org/10.1007/978-3-319-40229-1\\_8](http://dx.doi.org/10.1007/978-3-319-40229-1_8)
- [61] Nicolas Tabareau, Éric Tanter, and Matthieu Sozeau. 2017. Equivalences for Free! (July 2017). <https://hal.inria.fr/hal-01559073> working paper or preprint.
- [62] Iain Whiteside, David Aspinall, Lucas Dixon, and Gudmund Grov. 2011. Towards Formal Proof Script Refactoring. In

- Intelligent Computer Mathematics: 18th Symposium, Calculemus 2011, and 10th International Conference, MKM 2011, Bertinoro, Italy, July 18-23, 2011. Proceedings*, James H. Davenport, William M. Farmer, Josef Urban, and Florian Rabe (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 260–275. DOI:[http://dx.doi.org/10.1007/978-3-642-22673-1\\_18](http://dx.doi.org/10.1007/978-3-642-22673-1_18)
- [63] Thomas Williams, Pierre-Évariste Dagand, and Didier Rémy. 2014. Ornaments in Practice. In *Proceedings of the 10th ACM SIGPLAN Workshop on Generic Programming (WGP '14)*. ACM, New York, NY, USA, 15–24. DOI:<http://dx.doi.org/10.1145/2633628.2633631>
- [64] Doug Woos, James R. Wilcox, Steve Anton, Zachary Tatlock, Michael D. Ernst, and Thomas Anderson. 2016. 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)*. ACM, New York, NY, USA, 154–165. DOI:<http://dx.doi.org/10.1145/2854065.2854081>
- [65] Théo Zimmermann and Hugo Herbelin. 2015. Automatic and Transparent Transfer of Theorems along Isomorphisms in the Coq Proof Assistant. *CoRR* abs/1505.05028 (2015). <http://arxiv.org/abs/1505.05028>