

# Redirecting Proofs by Contradiction

Jasmin Christian Blanchette

Technische Universität München, Germany

## Abstract

This paper presents an algorithm that redirects proofs by contradiction. The input is a refutation graph, as produced by an automatic theorem prover (e.g., E, SPASS, Vampire, Z3); the output is a direct proof expressed in natural deduction extended with case analyses and nested subproofs. The algorithm is implemented in Isabelle’s Sledgehammer, where it enhances the legibility of machine-generated proofs.

## 1 Introduction

The proofs returned by automatic theorem provers (ATPs) are notoriously difficult to read. This is an issue if an ATP solves an open mathematical problem (such as the Robbins conjecture [11]), because users then certainly want to study the proof closely. But even in the context of program verification, where users are normally satisfied with a “proved” or “disproved,” they might still want to analyze proofs—for example, if they suspect errors in their axiom set.

Our interest in intelligible ATP proofs has a different origin. The tool Sledgehammer [19] integrates ATPs with the proof assistant Isabelle/HOL [15]. Given an Isabelle conjecture, Sledgehammer heuristically selects relevant lemmas from Isabelle’s libraries, translates them along with the conjecture to first-order logic, and sends the resulting problem to state-of-the-art provers such as E [24], SPASS [28], Vampire [22], and Z3 [5].

To guard against bugs in the ATPs and in Sledgehammer’s translation module, ATP proofs are reconstructed in Isabelle. This is accomplished through either a single invocation of the built-in resolution prover *metis* [8] or a structured Isar proof [20]. The latter option is useful for larger proofs, which *metis* fails to re-find within a reasonable time. But most users find the proofs unattractive and are disinclined to insert them in their theory text. As an illustration, consider the conjecture  $\text{length } (\text{tl } xs) \leq \text{length } xs$ , which states that the tail of a list (the list from which we remove its first element, or the empty list if the list is empty) is at most as long as the full list. The proof found by Vampire, translated to Isar, is as follows:

```
proof neg_clausify
  assume  $\text{length } (\text{tl } xs) \not\leq \text{length } xs$ 
  hence  $\text{drop } (\text{length } xs) (\text{tl } xs) \neq []$  by (metis drop_eq_Nil)
  hence  $\text{tl } (\text{drop } (\text{length } xs) xs) \neq []$  by (metis drop_tl)
  hence  $\forall u. xs @ u \neq xs \vee \text{tl } u \neq []$  by (metis append_eq_conv_conj)
  hence  $\text{tl } [] \neq []$  by (metis append_Nil2)
  thus False by (metis tl.simps(1))
qed
```

(The function `drop n` removes the first  $n$  elements from a list.) The `neg_clausify` proof method puts the Isabelle conjecture into negated clause form to ensure that it has the same shape as the corresponding ATP conjecture. The negation of the clause is introduced in the `assume` line, and a sequence of intermediate facts each introduced by `hence` leads to a contradiction.

There is a considerable body of research about making resolution proofs intelligible. Early work focused on translating detailed resolution proofs into natural deduction calculi [13, 21]. Although they are arguably more readable, these calculi operate at the logical level, whereas humans reason mostly at the “assertion level,” invoking definitions and lemmas without providing the full logical details. A line of research focused on transforming natural deduction proofs into assertion-level proofs [1, 7], culminating with the systems TRAMP [12] and Otterfier [31]. More related work includes the identification of obvious inferences [4, 23], the successful transformation of EQP’s proof of the Robbins conjecture using ILF [3], and more recently the use of TPTP-based tools to present Mizar articles [27].

It would have been interesting to try out TRAMP and Otterfier, but these are large pieces of unmaintained software that are hardly installable on modern machines and that only support older ATP systems. Regardless, the problem looks somewhat different in the context of Sledgehammer. Because the provers are given hundreds of lemmas as axioms, they tend to find short proofs with few lemmas. Moreover, Sledgehammer can coalesce consecutive inferences if short proofs are desired. Replaying an inference is a minor issue, thanks to *metis*.

The first obstacle to readability is that the Isar proof, like the underlying ATP proof, is by contradiction. This paper presents an algorithm for transforming proofs by contradiction into direct proofs—or *redirecting* proofs—to improve intelligibility. Knuth, Larrabee, and Roberts call the unnecessary use of proof by contradiction a sin against mathematical exposition [10, §3]—but since redirection is always possible, what would a necessary use look like?

The redirection algorithm is not tied to any one calculus or logic, as long as it admits contraposition. In particular, it works on the Isar proofs generated by Sledgehammer or directly on first-order TSTP proofs [26]. The direct proofs are expressed in a simple Isar-like syntax, which can be regarded as natural deduction extended with case analyses and nested subproofs (Section 2). The algorithm is first demonstrated on a few examples (Section 3) before it is presented in more detail, both in prose and as Standard ML pseudocode (Section 4).

For examples with a simple linear structure, such as the Isar proof above, the proof can be turned around by applying contraposition repeatedly:

```
proof -
  have tl [] = [] by (metis tl.simps(1))
  hence  $\exists u. xs @ u = xs \wedge tl\ u = []$  by (metis append_Nil2)
  hence tl (drop (length xs) xs) = [] by (metis append_eq_conv_conj)
  hence drop (length xs) (tl xs) = [] by (metis drop_tl)
  thus length (tl xs)  $\leq$  length xs by (metis drop_eq_Nil)
qed
```

The direct proof is easier to understand than the indirect one, even though it does not quite look like a human-written proof—humans would most likely avoid the detour through `drop`.

The approach works on arbitrary proofs by contradiction. A prototype demonstrated at earlier workshops [18, 19] sometimes exhibited exponential behavior. This has been resolved: Excluding a linear number of additional inferences that justify case analyses, each inference in the proof by contradiction now gives rise to exactly one inference in the direct proof. The algorithm can easily process proofs with hundreds or thousands of inferences.

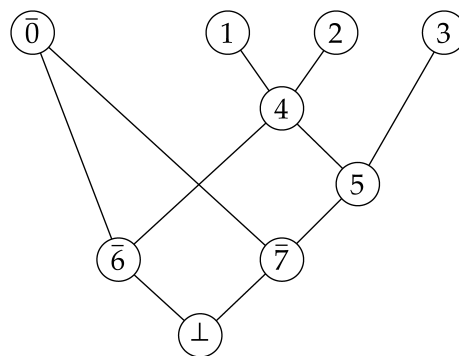
The algorithm is implemented in Sledgehammer. It is also described in Section 6.8 of my Ph.D. thesis [2], whose text largely forms the basis of this paper.

## 2 Proof Notations

**Proof Graph.** ATP proofs identify formulas by numbers. There may be several conjectures, in which case they are interpreted disjunctively. The negated conjectures and user-provided axioms are numbered  $0, 1, 2, \dots, n-1$ , and the derivations performed during proof search (whether or not they participate in the final proof) are numbered sequentially from  $n$ . We abstract the ATP proofs by ignoring the formulas and keeping only the numbers. We call formulas *atoms* since we are not interested in their structure. The letters  $a, b$  denote atoms.

An atom is *tainted* if it is one of the negated conjectures or has been derived, directly or indirectly, from a negated conjecture. For convenience, we relabel the ATP proof's atoms so that tainted atoms are decorated with a bar, denoting negation. Thus, if atom 3, corresponding to the formula  $\phi$ , is tainted, it is relabeled to  $\bar{3}$ , but it still stands for  $\phi$  and is called an atom despite the negative bar. After the relabeling, removing the bar negates the formula:  $\bar{3} \equiv \neg\phi$ .

A proof graph is a directed acyclic graph in which an edge  $a \rightarrow a'$  indicates that atom  $a$  is used to derive atom  $a'$ . Proof graphs are required to have exactly one sink node, whose formula is  $\perp$ , and only one connected component. We adopt the convention that derived nodes appear lower than their parent nodes in the graph and omit the arrowheads:



It is natural to write  $\perp$  rather than a numeric label for the sink node in examples.

In Sledgehammer, unary inferences are collapsed, and the first-order formulas are translated back to HOL before the proof is redirected. This is outside the scope of this paper and is explained in more detail in a companion paper [25].

**Isar Proofs.** Proof graphs cannot represent proofs by case analysis and only serve for the redirection algorithm’s input. We need more powerful notations for the output. Isar proofs [14, §4; 29] are a linear representation of natural deduction proofs in the style of Jaśkowski [9]. Unlike Gentzen-style trees [6], they allow the sharing of common derivations. The proof on the left-hand side is by contradiction; that on the right-hand side is the corresponding direct proof:

<pre> proof <i>neg_clausify</i>   assume <math>\bar{0}</math>   have 4 by (<i>metis</i> 1 2)   have 5 by (<i>metis</i> 3 4)   have <math>\bar{6}</math> by (<i>metis</i> <math>\bar{0}</math> 4)   have <math>\bar{7}</math> by (<i>metis</i> <math>\bar{0}</math> 5)   show <math>\perp</math> by (<i>metis</i> <math>\bar{6}</math> <math>\bar{7}</math>) qed                 </pre>	<pre> proof –   have 4 by (<i>metis</i> 1 2)   have 5 by (<i>metis</i> 3 4)   have <math>6 \vee 7</math> by <i>metis</i>   moreover   { assume 6     have 0 by (<i>metis</i> 4 6) }   moreover   { assume 7     have 0 by (<i>metis</i> 5 7) }   ultimately show 0 by <i>metis</i> qed                 </pre>
--	---

Notice that the direct proof involves a 2-way case analysis on a disjunction. Generalized disjunctions of the form  $a_1 \vee \dots \vee a_m$  are called *clauses* and are denoted by the letters  $c, d, e$ . Clauses are considered equal modulo associativity, commutativity, and idempotence. Sets of clauses are denoted by  $\Gamma$ .

Proof redirection requires that inferences can be redirected using the contrapositive but otherwise makes no assumptions about the proof calculus. Inferences that introduce new symbols can also be redirected; for example, skolemization becomes “un-herbrandization” [25, §4].

**Shorthand Proofs.** The last proof format is an ad hoc shorthand notation for a subset of Isar. In their simplest form, these shorthand proofs are a list of derivations  $c_1, \dots, c_m \triangleright c$  whose intuitive meaning is: “From the hypotheses  $c_1$  and ... and  $c_m$ , infer  $c$ .” The clauses on the left-hand side are interpreted as a set  $\Gamma$ .

If a hypothesis  $c_i$  is the previous derivation’s conclusion, we can omit it and write  $\blacktriangleright$  instead of  $\triangleright$ . This notation mimics Isar, with  $\triangleright$  for have or show and  $\blacktriangleright$  for hence or thus. Depending on whether we use the abbreviated format, our running example becomes

$1, 2 \triangleright 4$ $3, 4 \triangleright 5$ $\bar{0}, 4 \triangleright \bar{6}$ $\bar{0}, 5 \triangleright \bar{7}$ $\bar{6}, \bar{7} \triangleright \perp$	$1, 2 \triangleright 4$ $3 \blacktriangleright 5$ $\bar{0}, 4 \triangleright \bar{6}$ $\bar{0}, 5 \triangleright \bar{7}$ $\bar{6} \blacktriangleright \perp$
---	---

Each derivation  $\Gamma \triangleright c$  is essentially a sequent with  $\Gamma$  as the antecedent and  $c$  as the succedent. For proofs by contradiction, the clauses in the antecedent are either the negated conjecture ( $\bar{0}$ ),

atoms that correspond to background facts (1, 2, and 3), or atoms that were proved in preceding sequents (4, 5,  $\bar{6}$ , and  $\bar{7}$ ); the succedent of the last sequent is always  $\perp$ .

Direct proofs can be presented in the same way, but the negated conjecture  $\bar{0}$  may not appear in any of the sequents' antecedents, and the last sequent must have the conjecture 0 as its succedent. In some of the direct proofs, it is useful to introduce case analyses. For example:

$$\begin{array}{c} 1, 2 \triangleright 4 \\ 3 \blacktriangleright 5 \\ \triangleright 6 \vee 7 \\ \left[ \begin{array}{c|c} [6] & [7] \\ 4 \blacktriangleright 0 & 5 \blacktriangleright 0 \end{array} \right] \end{array}$$

In general, case analysis blocks have the form

$$\left[ \begin{array}{c|c|c} [c_1] & \dots & [c_m] \\ \Gamma_{11} \triangleright d_{11} & \dots & \Gamma_{m1} \triangleright d_{m1} \\ \vdots & & \vdots \\ \Gamma_{1k_1} \triangleright d_{1k_1} & \dots & \Gamma_{mk_m} \triangleright d_{mk_m} \end{array} \right]$$

with the requirement that a sequent with the succedent  $c_1 \vee \dots \vee c_m$  has been proved immediately above the case analysis. Each of the branches must also be a valid proof. The assumptions  $[c_i]$  may be used to discharge hypotheses in the same branch, as if they had been sequents  $\triangleright c_i$ . The case analysis will sometimes be regarded as a sequent

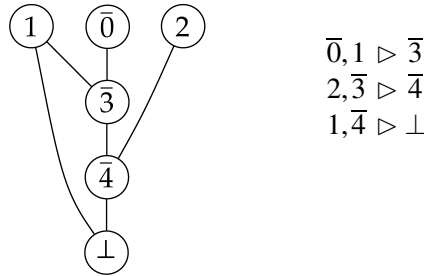
$$c_1 \vee \dots \vee c_m, \bigcup_{i,j} (\Gamma_{ij} - c_i - \bigcup_{j' < j} d_{ij'}) \triangleright d_{1k_1} \vee \dots \vee d_{mk_m}$$

by ignoring its internal structure.

### 3 Examples of Proof Redirection

Before reviewing the redirection algorithm, we consider four examples of proofs by contradiction and redirect them to produce a direct proof. The first example has a simple linear structure, the second and third examples involve a ‘‘lasso,’’ and the last example has no apparent structure.

**A Linear Proof.** We start with a simple proof by contradiction expressed as a proof graph and in our shorthand notation:



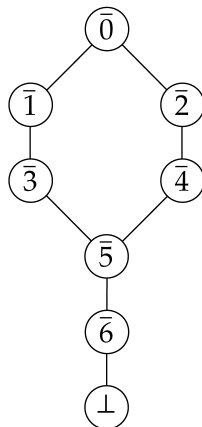
We redirect the sequents using sequent-level contraposition to eliminate all taints (represented as bars after the relabeling). This gives

$$\begin{aligned} 1, 3 &\triangleright 0 \\ 2, 4 &\triangleright 3 \\ 1 &\triangleright 4 \end{aligned}$$

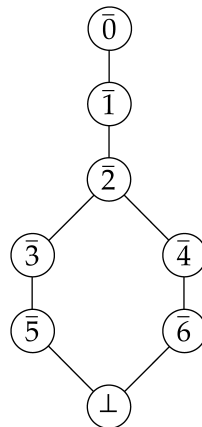
We then obtain the direct proof by reversing the order of the sequents, and introduce  $\blacktriangleright$  wherever possible:

proof –  
 have 4 by (*metis* 1)             $1 \triangleright 4$   
 hence 3 by (*metis* 2)            $2 \blacktriangleright 3$   
 thus 0 by (*metis* 1)             $1 \blacktriangleright 0$   
 qed

**Lasso-Shaped Proofs.** The next two examples look superficially like lassos but are of course acyclic, as required of all proof graphs:



$$\begin{aligned} \bar{0} &\triangleright \bar{1} \\ \bar{0} &\triangleright \bar{2} \\ \bar{1} &\triangleright \bar{3} \\ \bar{2} &\triangleright \bar{4} \\ \bar{3}, \bar{4} &\triangleright \bar{5} \\ \bar{5} &\triangleright \bar{6} \\ \bar{6} &\triangleright \perp \end{aligned}$$



$$\begin{aligned} \bar{0} &\triangleright \bar{1} \\ \bar{1} &\triangleright \bar{2} \\ \bar{2} &\triangleright \bar{3} \\ \bar{2} &\triangleright \bar{4} \\ \bar{3} &\triangleright \bar{5} \\ \bar{4} &\triangleright \bar{6} \\ \bar{5}, \bar{6} &\triangleright \perp \end{aligned}$$

We start with the example on the left-hand side. Starting from  $\perp$ , it is easy to redirect the stem:

$$\begin{aligned} &\triangleright 6 \\ 6 &\triangleright 5 \\ 5 &\triangleright 3 \vee 4 \end{aligned}$$

When applying the contrapositive to eliminate the negations in  $\bar{3}, \bar{4} \triangleright \bar{5}$ , we obtain a disjunction in the succedent:  $5 \triangleright 3 \vee 4$ . To continue from there, we introduce a case analysis. In each branch, we can finish the proof:

$$\left[ \begin{array}{c|c} [3] & [4] \\ 3 \triangleright 1 & 4 \triangleright 2 \\ 1 \triangleright 0 & 2 \triangleright 0 \end{array} \right]$$

In the second lasso example, the cycle occurs near the end of the contradiction proof. A disjunction already arises when we redirect the last derivation. Naively finishing each branch independently leads to a fair amount of duplication:

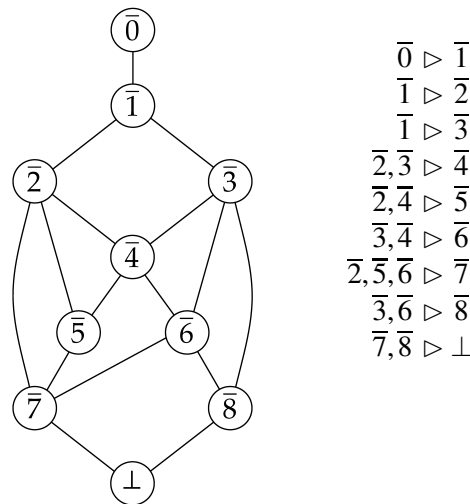
$$\begin{array}{c} \triangleright 5 \vee 6 \\ \left[ \begin{array}{c|c} [5] & [6] \\ 5 \triangleright 3 & 6 \triangleright 4 \\ 3 \triangleright 2 & 4 \triangleright 2 \\ 2 \triangleright 1 & 2 \triangleright 1 \\ 1 \triangleright 0 & 1 \triangleright 0 \end{array} \right] \end{array}$$

The key observation is that the two branches can share the last two inferences. This yields the following proof (without and with  $\blacktriangleright$ ):

$$\begin{array}{c} \triangleright 5 \vee 6 \\ \left[ \begin{array}{c|c} [5] & [6] \\ 5 \triangleright 3 & 6 \triangleright 4 \\ 3 \triangleright 2 & 4 \triangleright 2 \\ 2 \triangleright 1 \\ 1 \triangleright 0 \end{array} \right] \end{array} \qquad \begin{array}{c} \triangleright 5 \vee 6 \\ \left[ \begin{array}{c|c} [5] & [6] \\ \blacktriangleright 3 & \blacktriangleright 4 \\ \blacktriangleright 2 & \blacktriangleright 2 \\ \blacktriangleright 1 \\ \blacktriangleright 0 \end{array} \right] \end{array}$$

Here we were fortunate that the branches were joinable on the atom 2. To avoid duplication, we must in general join on a disjunction  $a_1 \vee \dots \vee a_m$ , as in the next example.

**A Spaghetti Proof.** The final example is diabolical (and slightly unrealistic, perhaps):



We start with the contrapositive of the last sequent:

$$\triangleright 7 \vee 8$$

We perform a case analysis on  $7 \vee 8$ . Since we want to avoid duplication in the two branches, we first determine which nodes are reachable in the refutation graph by navigating upward from either  $\bar{7}$  or  $\bar{8}$  but not from both. The only such nodes here are  $\bar{5}$ ,  $\bar{7}$ , and  $\bar{8}$ . In each branch, we can perform derivations of the form  $\Gamma \triangleright b$  where  $\Gamma \cap \{5, 7, 8\} \neq \emptyset$  without fearing duplication. Following this rule, we can only perform one inference in the right branch before we must stop:

$$\begin{array}{c} [8] \\ 8 \triangleright 3 \vee 6 \end{array}$$

Any further inferences would need to be repeated in the left branch, so it is indeed a good idea to stop. The left branch starts as follows:

$$\begin{array}{c} [7] \\ 7 \triangleright 2 \vee 5 \vee 6 \end{array}$$

We would now like to perform the inference  $5 \triangleright 2 \vee 4$ . This would certainly not lead to any duplication, because  $\bar{5}$  is not reachable from  $\bar{8}$  by navigating upward in the refutation graph. However, we cannot discharge the hypothesis 5, having established only the disjunction  $2 \vee 5 \vee 6$ . We need a case analysis on the disjunction to proceed:

$$\left[ \begin{array}{c|c} [2] & \begin{array}{c} [5] \\ 5 \triangleright 2 \vee 4 \end{array} \\ \hline & [6] \end{array} \right]$$

The 2 and 6 subbranches are left alone, because there is no node that is reachable only from  $\bar{2}$  or  $\bar{6}$  but not from the other two nodes in  $\{\bar{2}, \bar{5}, \bar{6}\}$  by navigating upward in the refutation graph. Since only one branch is nontrivial, it is arguably more aesthetically pleasing to abbreviate the entire case analysis to

$$2 \vee 5 \vee 6 \triangleright 2 \vee 4 \vee 6$$

Putting this all together, the outer case analysis becomes

$$\left[ \begin{array}{c|c} [7] & [8] \\ \hline \blacktriangleright 2 \vee 5 \vee 6 & \blacktriangleright 3 \vee 6 \\ \blacktriangleright 2 \vee 4 \vee 6 & \end{array} \right]$$

The left branch proves  $2 \vee 4 \vee 6$ , the right branch proves  $3 \vee 6$ ; hence, both branches together prove  $2 \vee 3 \vee 4 \vee 6$ . Next, we perform the inference  $6 \triangleright 3 \vee 4$ . This requires a case analysis on  $2 \vee 3 \vee 4 \vee 6$ :

$$\left[ \begin{array}{c|c|c|c} [2] & [3] & [4] & \begin{array}{c} [6] \\ 6 \triangleright 3 \vee 4 \end{array} \\ \hline & & & \end{array} \right]$$

This proves  $2 \vee 3 \vee 4$ . Since only one branch is nontrivial, we prefer to abbreviate the case analysis to

$$2 \vee 3 \vee 4 \vee 6 \triangleright 2 \vee 3 \vee 4$$



It may help to think of such abbreviated inferences as instances of rewriting modulo associativity, commutativity, and idempotence. Here, 6 is rewritten to  $3 \vee 4$  in  $2 \vee 3 \vee 4 \vee 6$ , resulting in  $2 \vee 3 \vee 4$ . Similarly, the sequent  $4 \triangleright 2 \vee 3$  gives rise to the case analysis

$$\left[ \begin{array}{c|c} [2] & [3] \\ \hline 4 \triangleright 2 \vee 3 & [4] \end{array} \right]$$

which can be abbreviated as well. We are left with  $2 \vee 3$ . The rest is analogous to the second lasso-shaped proof:

$$\left[ \begin{array}{c|c} [2] & [3] \\ \hline 2 \triangleright 1 & 3 \triangleright 1 \end{array} \right] \\ 1 \triangleright 0$$

Putting all of this together, we obtain the following proof, expressed in Isar and in shorthand. The result is quite respectable, considering the spaghetti-like graph we started with:

```

proof -
  have 7 ∨ 8 by metis
  moreover
  { assume 7
    hence 2 ∨ 5 ∨ 6 by metis
    hence 2 ∨ 4 ∨ 6 by metis }
  moreover
  { assume 8
    hence 3 ∨ 6 by metis }
  ultimately have 2 ∨ 3 ∨ 4 ∨ 6 by metis
  hence 2 ∨ 3 ∨ 4 by metis
  hence 2 ∨ 3 by metis
  moreover
  { assume 2
    hence 1 by metis }
  moreover
  { assume 3
    hence 1 by metis }
  ultimately have 1 by metis
  thus 0 by metis
qed
    
```

$$\begin{array}{c} \triangleright 7 \vee 8 \\ \left[ \begin{array}{c|c} [7] & [8] \\ \hline \blacktriangleright 2 \vee 5 \vee 6 & \blacktriangleright 3 \vee 6 \\ \blacktriangleright 2 \vee 4 \vee 6 & \blacktriangleright 2 \vee 3 \vee 4 \\ & \blacktriangleright 2 \vee 3 \\ & \blacktriangleright 0 \end{array} \right] \\ \left[ \begin{array}{c|c} [2] & [3] \\ \hline \blacktriangleright 1 & \blacktriangleright 1 \end{array} \right] \\ \blacktriangleright 0 \end{array}$$

## 4 The Redirection Algorithm

The process we applied in the examples above can be generalized into an algorithm. The algorithm takes an arbitrary proof by contradiction expressed as a set of sequents as input, and produces a proof in our Isar-like shorthand notation, with sequents and case analysis blocks. The proof is constructed one inference at a time starting from  $\top$  (the negation of  $\perp$ ) until the conjecture—in general, the disjunction of the conjectures—is proved.

**Basic Concepts.** A fundamental operation is sequent-level contraposition. Let  $a_1, \dots, a_m$  be the untainted atoms and  $\bar{b}_1, \dots, \bar{b}_n$  the tainted atoms of a proof by contradiction. The proof then consists of the following three kinds of sequent (with  $n > 0$ ):

$$a_1, \dots, a_m, \bar{b}_1, \dots, \bar{b}_n \triangleright \perp \quad a_1, \dots, a_m, \bar{b}_1, \dots, \bar{b}_n \triangleright \bar{b} \quad a_1, \dots, a_m \triangleright a$$

Their *contrapositives* are, respectively,

$$a_1, \dots, a_m \triangleright b_1 \vee \dots \vee b_n \quad a_1, \dots, a_m, b \triangleright b_1 \vee \dots \vee b_n \quad a_1, \dots, a_m \triangleright a$$

We call the contrapositives of the sequents in the proof by contradiction the *redirected sequents*.

Based on the set of redirected sequents, we define the *atomic inference graph* (AIG) with, for each redirected sequent  $\Gamma \triangleright c$ , an edge from each atom in  $\Gamma$  to each atom in  $c$ , and no additional edges. The AIG encodes the order in which the atoms can be inferred in a direct proof. Navigating forward (downward) in this graph along the unnegated tainted atoms  $b_j$  corresponds to navigating backward (upward) in the refutation graph along the  $\bar{b}_j$ 's.

Like the underlying refutation graph, the AIG is acyclic and connected. Potential cycles would involve either only untainted atoms  $a_i$ , only tainted atoms  $b_j$ 's, or a mixture of both kinds. A cycle  $a_{i_1} \rightarrow \dots \rightarrow a_{i_k} \rightarrow a_{i_1}$  is impossible, because the contrapositive leaves these inferences unchanged and hence the cycle would need to occur in the (acyclic) refutation graph. A cycle  $b_{j_1} \rightarrow \dots \rightarrow b_{j_k} \rightarrow b_{j_1}$  is impossible, because the contrapositive turns all the edges around and hence the reverse cycle would need to occur in the refutation graph. Finally, mixed cycles necessarily involve an edge  $b \rightarrow a$ , which is impossible because redirected sequents with untainted atoms  $a$  can only have untainted atoms as predecessors (cf.  $a_1, \dots, a_m \triangleright a$ ).

Given a set of (tainted or untainted) atoms  $A$ , the *zone* of an atom  $a \in A$  with respect to  $A$  is the set of possibly trivial descendants of  $a$  in the AIG that are not descendants of any of the other atoms in  $A$ . As a trivial descendant of itself,  $a$  will either belong to its own zone or to no zone at all (depending on whether it is a descendant of a node  $a' \in A - \{a\}$ ). Zones identify inferences that can safely be performed inside a branch in a case analysis.

**The Algorithm.** The algorithm keeps track of the *last-proved clause* (initially  $\top$ ), the set of *already proved atoms* (initially the set of facts taken as axioms), and the set of *remaining sequents* to use (initially all the redirected sequents provided as input). It performs these steps:

1. If there are no remaining sequents, stop.
2. If the last-proved clause is  $\top$  or a single atom:
  - 2.1. Pick a sequent  $\Gamma \triangleright c$  among the remaining sequents that can be proved using only already proved atoms, preferring sequents with a single atom in their succedent.
  - 2.2. Append  $\Gamma \triangleright c$  to the proof.
  - 2.3. Make  $c$  the last-proved clause, add  $c$  to the already proved atoms if it is an atom, and remove  $\Gamma \triangleright c$  from the remaining sequents.
  - 2.4. Go to step 1.

3. Otherwise, the last-proved succedent is of the form  $a_1 \vee \dots \vee a_m$ . An  $m$ -way case analysis is called for:<sup>1</sup>
  - 3.1. Compute the zone of each atom  $a_i$  with respect to  $\{a_1, \dots, a_m\}$ .
  - 3.2. For each  $a_i$ , compute the set  $\mathcal{S}_i$  of sequents  $\Gamma \triangleright c$  such that  $\Gamma$  consists only of already proved atoms or atoms within  $a_i$ 's zone.
  - 3.3. Recursively invoke the algorithm  $m$  times, once for each  $a_i$ , each time with  $a_i$  as the last-proved clause,  $a_i$  added to the already proved atoms, and  $\mathcal{S}_i$  as the set of remaining sequents. This step yields  $m$  (possibly empty) subproofs  $\pi_1, \dots, \pi_m$ .
  - 3.4. Append the following case analysis block to the proof:
 
$$\left[ \begin{array}{c|c|c} [a_1] & \dots & [a_m] \\ \pi_1 & \dots & \pi_m \end{array} \right]$$
  - 3.5. Make the succedent  $b_1 \vee \dots \vee b_n$  of the case analysis block (regarded as a sequent) the last-proved clause, add  $b_1$  to the already proved atoms if  $k = 1$ , and remove all sequents belonging to any of the sets  $\mathcal{S}_i$  from the remaining sequents.
  - 3.6. Go to step 1.

Whenever a redirected sequent is generated, it is removed from the set of remaining sequents. In step 3, the recursive calls operate on pairwise disjoint subsets  $\mathcal{S}_i$  of the remaining sequents. Consequently, each redirected sequent appears at most once in the generated proof, and the resulting direct proof contains the same number of inferences as the initial proof by contradiction. In Isar, each case analysis is additionally justified by one *metis* inference.

In the degenerate case where no atoms are tainted (i.e., the proof exploits an inconsistency in the axiom set), the generated proof is simply a linearization of the refutation graph, and the last inference proves  $\perp$  (which is, unusually, untainted). To produce a syntactically valid Isar proof, a trivial inference must be added to derive the conjecture from  $\perp$ .

**Pseudocode.** To make the above description more concrete, the algorithm is presented in Standard ML pseudocode below.<sup>2</sup> The pseudocode is fairly faithful to the description above. Atoms are represented by integers and literals by sets (lists) of integers. Go-to statements are implemented by recursion, and the state is threaded through recursive calls as three arguments (*last*, *earlier*, and *seqs*). One notable difference, justified by a desire to avoid code duplication, is that the set of already proved atoms, called *earlier*, excludes the last-proved clause *last*. Hence, we take  $last \cup earlier$  to obtain the already proved atoms, where *last* is either the empty list (representing  $\top$ ) or a singleton list (representing a single atom).

Shorthand proofs are represented as lists of *inferences*:

`datatype inference = Have of int list  $\times$  int list | Cases of (int  $\times$  inference list) list`

<sup>1</sup>A generalization would be to perform a  $m'$ -way case analysis, with  $m' < m$ , by keeping some disjunctions. For example, we could perform a 3-way case analysis with  $a_1 \vee a_2$ ,  $a_3$ , and  $a_4$  as the assumptions instead of breaking all the disjunctions in a 4-way analysis. This could lead to nicer proofs if the disjuncts are carefully chosen.

<sup>2</sup>The actual ML code is distributed with Isabelle. A recent repository version is available at [http://isabelle.in.tum.de/repos/isabelle/file/e5303bd748f2/src/HOL/Tools/ATP/atp\\_proof\\_redirect.ML](http://isabelle.in.tum.de/repos/isabelle/file/e5303bd748f2/src/HOL/Tools/ATP/atp_proof_redirect.ML).

The main function implementing the algorithm follows:

```

fun redirect last earlier seqs =
  if null seqs then
    []
  else if length last ≤ 1 then
    let val provable = filter (fn (Γ, _) ⇒ Γ ⊆ last ∪ earlier) seqs
        val horn_provable = filter (fn (_, [c]) ⇒ true | _ ⇒ false) provable
        val (Γ, c) = hd (horn_provable @ provable)
    in Have (Γ, c) :: redirect c (last ∪ earlier) (seqs - {(Γ, c)}) end
  else
    let val zs = zones_of (length last) (map (descendants seqs) last)
        val ℳ = map (fn z ⇒ filter (fn (Γ, _) ⇒ Γ ⊆ earlier ∪ z) seqs) zs
        val cases = map (fn (a, ss) ⇒ (a, redirect [a] earlier ss)) (zip last ℳ)
    in Cases cases :: redirect (succedent_of_cases cases) earlier (seqs - ⋃ ℳ) end

```

The code uses familiar ML functions, such as map, filter, and zip. It also relies on a descendants function that returns the descendants of the specified node in the AIG associated with *seqs*; its definition is omitted. Finally, the code depends on the following straightforward functions:

```

fun zones_of 0 _ = []
  | zones_of n (B :: Bs) = (B - ⋃ Bs) :: zones_of (n - 1) (Bs @ [B])

fun succedent_of_inf (Have (_, c)) = c
  | succedent_of_inf (Cases cases) = succedent_of_cases cases
and succedent_of_case (a, []) = [a]
  | succedent_of_case (_, infs) = succedent_of_inf (last infs)
and succedent_of_cases cases = ⋃ (map succedent_of_case cases)

```

**Correctness.** It is not hard to convince ourselves that the proof output by `redirect` is correct by inspecting the code. A `Have (Γ, c)` sequent is appended only if all the atoms in  $\Gamma$  have been proved (or assumed) already, and a case analysis on  $a_1 \vee \dots \vee a_m$  always follows a sequent with the succedent  $a_1 \vee \dots \vee a_m$ . Whenever a sequent is output, it is removed from *seqs*. The function returns only if *seqs* is empty, at which point the conjecture must have been proved (except in the degenerate case where the negated conjecture does not participate in the refutation).

Termination is not quite as obvious. The recursion is well-founded, because the pair  $(\text{length } seqs, \text{length } last)$  becomes strictly smaller with respect to the lexicographic extension of  $<$  on natural numbers for each of the three syntactic recursive calls.

- For the first recursive call, the list  $seqs - \{(\Gamma, c)\}$  is strictly shorter than *seqs* since  $(\Gamma, c) \in seqs$ .
- The second call is performed for each branch of a case analysis; the *ss* argument is a (not necessarily strict) subset of the caller's *seqs*, and the list  $[a]$  is strictly shorter than *last*, which has length 2 or more.

- For the third call, the key property is that at least one of the zones is nonempty, from which we obtain  $seqs - \bigcup \mathcal{S} \subset seqs$ . If all the zones were empty, each atom  $a_i$  would be the descendant of at least one atom  $a_{i'}$  in the AIG (with  $i' \neq i$ ), which is impossible because the AIG is acyclic.

As for run-time exceptions, the only worrisome construct is the hd call in `redirect`'s second branch. We must convince ourselves that there exists at least one sequent  $(\Gamma, c) \in seqs$  such that  $\Gamma \subseteq last \cup earlier$ . Intuitively, this is unsurprising, because  $seqs$  is initialized from a well-formed refutation graph: The nonexistence of such a sequent would indicate a gap or a cycle in the refutation graph. More precisely, if there exist untainted atoms  $\notin last \cup earlier$ , these can always be processed first; indeed, the preference for sequents with a single atom in their succedent ensures that they are processed before the first case analysis. Otherwise:

- If  $last$  is  $\square$  (representing  $\top$ ) or an untainted atom, the contrapositive  $a_1, \dots, a_m \triangleright b_1 \vee \dots \vee b_n$  of the very last inference is applicable since it only depends on untainted atoms, all of which have already been proved.
- Otherwise,  $last$  is a tainted atom  $b$ . The refutation graph must contain an inference  $a_1, \dots, a_m, \bar{b}_1, \dots, \bar{b}_n \triangleright \bar{b}$ , whose redirected inference is  $a_1, \dots, a_m, b \triangleright b_1 \vee \dots \vee b_n$ . Since it only depends on  $b$  and untainted atoms, it is applicable.

**Inlining.** As a postprocessing step, we can abbreviate case analyses in which only one branch is nontrivial, transforming

$$\left[ \begin{array}{c|c|c|c|c|c|c} & & & & [c_i] & & \\ & & & d_{11}, \dots, d_{1k_1} \triangleright e_1 & & & \\ & & & \vdots & & & \\ [c_1] & \cdots & [c_{i-1}] & d_{n1}, \dots, d_{nk_n} \triangleright e_n & [c_{i+1}] & \cdots & [c_m] \end{array} \right] \text{ into } \begin{array}{l} \widetilde{d}_{11}, \dots, \widetilde{d}_{1k_1} \triangleright \widetilde{e}_1 \\ \vdots \\ \widetilde{d}_{n1}, \dots, \widetilde{d}_{nk_n} \triangleright \widetilde{e}_n \end{array}$$

where the function  $\sim$  is the identity except for the assumption  $c_i$  and the conclusions  $e_1, \dots, e_n$ :

$$\widetilde{c}_i = c_1 \vee \dots \vee c_m \quad \widetilde{e}_j = c_1 \vee \dots \vee c_{i-1} \vee e_j \vee c_{i+1} \vee \dots \vee c_m$$

It is debatable whether such inlining is a good idea. The resulting proof has a simpler structure, with fewer nested proof blocks. However, these nested blocks can help make complex proof more intelligible. Moreover, the  $n$ -fold repetition of the disjuncts  $c_1, \dots, c_{i-1}, c_{i+1}, \dots, c_m$  clutters the proof and can slow it down.

The above procedure can be generalized to arbitrary case analysis blocks. I am grateful to Konstantin Korovin for insights that lead me to realize this. We can rewrite an  $m$ -way case analysis

$$\left[ \begin{array}{c|c|c} [c_1] & \cdots & [c_m] \\ \Gamma_{11} \triangleright e_{11} & \cdots & \Gamma_{m1} \triangleright e_{m1} \\ \vdots & & \vdots \\ \Gamma_{1k_1} \triangleright e_{1k_1} & \cdots & \Gamma_{mk_m} \triangleright e_{mk_m} \end{array} \right]$$

into a sequence of  $m$  case analyses:

$$\left[ \begin{array}{c|c|c|c|c|c|c} & & & & [c_i] & & \\ & & & \Gamma_{i1} \triangleright e_{i1} & & & \\ & & & \vdots & & & \\ [e_{1k_1}] & \cdots & [e_{(i-1)k_{i-1}}] & \Gamma_{ik_i} \triangleright e_{ik_i} & [c_{i+1}] & \cdots & [c_m] \end{array} \right]$$

Each of these has only one nontrivial branch and can be inlined, yielding a branch-free proof. For the spaghetti proof of the previous section, this process yields

$$\begin{aligned} &\triangleright 7 \vee 8 \\ &\blacktriangleright 2 \vee 5 \vee 6 \vee 8 \\ &\blacktriangleright 2 \vee 4 \vee 6 \vee 8 \\ &\blacktriangleright 2 \vee 3 \vee 4 \vee 6 \\ &\blacktriangleright 2 \vee 3 \vee 4 \\ &\blacktriangleright 2 \vee 3 \\ &\blacktriangleright 1 \vee 3 \\ &\blacktriangleright 1 \\ &\blacktriangleright 0 \end{aligned}$$

The example shows clearly that we rapidly obtain large disjunctions. In practice, each of the disjuncts would be an arbitrarily complex formula. Local definitions could be used to avoid repeating the formulas, but the loss of modularity is deplorable. Indeed, similar concerns about Hoare-style proof outlines for separation logic have led to the development of ribbon proofs [30], whose parallel ribbons evoke the branches of a case analysis.

If branch-free proofs are nonetheless desired, they can be generated more directly by iteratively “rewriting” the atoms, following a suggestion by Korovin. For example, starting from the sequent  $\triangleright 7 \vee 8$ , rewriting 7 would involve resolving  $\triangleright 7 \vee 8$  with  $7 \triangleright 2 \vee 5 \vee 6$ , resulting in  $2 \vee 5 \vee 6 \vee 8$ . In general, rewriting a tainted atom  $b_j$  within a sequent  $\Gamma \triangleright b_1 \vee \cdots \vee b_n$  involves resolving that sequent with the redirected sequent that has  $b_j$  in its assumptions. To guarantee that the procedure is linear, it suffices to rewrite atoms only if all their ancestors in the AIG have already been rewritten, thereby preventing multiple rewrites of the same atom.

## 5 Conclusion

This paper presented an algorithm that transforms proofs by contradiction as returned by automatic theorem provers into direct proofs. It sometimes introduces case splits but avoids duplicating inferences in the different branches of the split by joining again as early as possible. The resulting proofs are direct Isar proofs that have some of the structure Isabelle users have come to expect. The described procedure is admittedly fairly straightforward; it would not be surprising if it were part of folklore or a special case of existing work.

While the output is designed for replaying proofs, it also has a pedagogical value: Unlike Isabelle’s automatic tactics, which are black boxes, the proofs delivered by Sledgehammer can

be inspected and understood. The direct proof also forms a good basis for manual tuning. Further steps toward robust, intelligible Isar proofs, are described in a companion paper [25]. In future work, I am interested in transformations that increase proof readability and in the automatic discovery of concepts and lemmas, such as those available for Mizar proofs [16, 17].

**Acknowledgment.** Tobias Nipkow, Lawrence Paulson, Geoff Sutcliffe, and Josef Urban motivated me to pursue this research. Andrei Popescu helped simplify the approach considerably. Charles Francis, sponsored by Google Summer of Code, developed a prototype. Sascha Böhme, Konstantin Korovin, Tobias Nipkow, Hans de Nivelle, Josef Urban, Tjark Weber, and two (still) anonymous reviewers produced a wealth of comments that helped enhance this text. I thank them all. The research was financed by the Deutsche Forschungsgemeinschaft (DFG) projects Quis Custodiet (grant Ni 491/11-1) and Hardening the Hammer (grant Ni 491/14-1).

## References

- [1] S. Autexier, C. Benzmüller, A. Fiedler, H. Horacek, and Q. B. Vo. Assertion-level proof representation with under-specification. *Electr. Notes Theor. Comput. Sci.*, 93:5–23, 2004.
- [2] J. C. Blanchette. *Automatic Proofs and Refutations for Higher-Order Logic*. Ph.D. thesis, Dept. of Informatics, Technische Universität München, 2012.
- [3] B. I. Dahn. Robbins algebras are Boolean: A revision of McCune’s computer-generated solution of Robbins problem. *J. Algebra*, 208(2):526–532, 1998.
- [4] M. Davis. Obvious logical inferences. In P. J. Hayes, editor, *IJCAI ’81*, pages 530–531. William Kaufmann, 1981.
- [5] L. M. de Moura and N. Bjørner. Z3: An efficient SMT solver. In C. R. Ramakrishnan and J. Rehof, editors, *TACAS 2008*, volume 4963 of *LNCS*, pages 337–340. Springer, 2008.
- [6] G. Gentzen. Untersuchungen über das logische Schließen. *Math. Z.*, 39:176–210, 1935.
- [7] X. Huang. Translating machine-generated resolution proofs into ND-proofs at the assertion level. In N. Y. Foo and R. Goebel, editors, *PRICAI ’96*, volume 1114 of *LNCS*, pages 399–410. Springer, 1996.
- [8] J. Hurd. First-order proof tactics in higher-order logic theorem provers. In M. Archer, B. Di Vito, and C. Muñoz, editors, *Design and Application of Strategies/Tactics in Higher Order Logics*, number CP-2003-212448 in NASA Technical Reports, pages 56–68, 2003.
- [9] S. Jaśkowski. On the rules of suppositions in formal logic. *Studia Logica*, 1:5–32, 1934.
- [10] D. E. Knuth, T. L. Larrabee, and P. M. Roberts. *Mathematical Writing*. Mathematical Association of America, 1989.
- [11] W. McCune. Solution of the Robbins problem. *J. Autom. Reasoning*, 19(3):263–276, 1997.
- [12] A. Meier. TRAMP: Transformation of machine-found proofs into natural deduction proofs at the assertion level (system description). In D. McAllester, editor, *CADE-17*, volume 1831 of *LNAI*, pages 460–464. Springer, 2000.
- [13] D. Miller and A. Felty. An integration of resolution and natural deduction theorem proving. In *AAAI-86*, volume I: Science, pages 198–202. Morgan Kaufmann, 1986.
- [14] T. Nipkow. Programming and proving in Isabelle. <http://isabelle.in.tum.de/dist/Isabelle2013/doc/prog-prove.pdf>, 2013.

- [15] T. Nipkow, L. C. Paulson, and M. Wenzel. *Isabelle/HOL: A Proof Assistant for Higher-Order Logic*, volume 2283 of *LNCS*. Springer, 2002.
- [16] K. Pąk. The methods of improving and reorganizing natural deduction proofs. In *MathUI10*, 2010.
- [17] K. Pąk. Methods of lemma extraction in natural deduction proofs. *J. Autom. Reasoning*, 50(2):217–228, 2013.
- [18] L. C. Paulson. Three years of experience with Sledgehammer, a practical link between automated and interactive theorem provers. In B. Konev, R. Schmidt, and S. Schulz, editors, *PAAR-2010*, pages 1–10, 2010.
- [19] L. C. Paulson and J. C. Blanchette. Three years of experience with Sledgehammer, a practical link between automatic and interactive theorem provers. In G. Sutcliffe, E. Ternovska, and S. Schulz, editors, *IWIL-2010*, 2010.
- [20] L. C. Paulson and K. W. Susanto. Source-level proof reconstruction for interactive theorem proving. In K. Schneider and J. Brandt, editors, *TPHOLs 2007*, volume 4732 of *LNCS*, pages 232–245. Springer, 2007.
- [21] F. Pfenning. Analytic and non-analytic proofs. In R. E. Shostak, editor, *CADE-7*, volume 170 of *LNCS*, pages 393–413. Springer, 1984.
- [22] A. Riazanov and A. Voronkov. The design and implementation of Vampire. *AI Comm.*, 15(2-3):91–110, 2002.
- [23] P. Rudnicki. Obvious inferences. *J. Autom. Reasoning*, 3(4):383–393, 1987.
- [24] S. Schulz. System description: E 0.81. In D. Basin and M. Rusinowitch, editors, *IJCAR 2004*, volume 3097 of *LNAI*, pages 223–228. Springer, 2004.
- [25] S. J. Smolka and J. C. Blanchette. Robust, semi-intelligible Isabelle proofs from ATP proofs. In J. C. Blanchette and J. Urban, editors, *PxTP 2013*, 2013.
- [26] G. Sutcliffe, J. Zimmer, and S. Schulz. TSTP data-exchange formats for automated theorem proving tools. In W. Zhang and V. Sorge, editors, *Distributed Constraint Problem Solving and Reasoning in Multi-Agent Systems*, volume 112 of *Frontiers in Artificial Intelligence and Applications*, pages 201–215. IOS Press, 2004.
- [27] J. Urban, G. Sutcliffe, S. Trac, and Y. Puzis. Combining Mizar and TPTP semantic presentation and verification tools. In A. Grabowski and A. Naumowicz, editors, *Computer Reconstruction of the Body of Mathematics*, volume 18(31) of *Studies in Logic, Grammar and Rhetoric*, pages 121–136. University of Białystok, 2009.
- [28] C. Weidenbach. Combining superposition, sorts and splitting. In A. Robinson and A. Voronkov, editors, *Handbook of Automated Reasoning*, volume II, pages 1965–2013. Elsevier, 2001.
- [29] M. Wenzel. Isabelle/Isar—A generic framework for human-readable proof documents. In R. Matuszewski and A. Zalewska, editors, *From Insight to Proof—Festschrift in Honour of Andrzej Trybulec*, volume 10(23) of *Studies in Logic, Grammar and Rhetoric*. University of Białystok, 2007.
- [30] J. Wickerson, M. Dodds, and M. J. Parkinson. Ribbon proofs for separation logic. In M. Felleisen and P. Gardner, editors, *ESOP 2013*, volume 7792 of *LNCS*, pages 189–208. Springer, 2013.
- [31] J. Zimmer, A. Meier, G. Sutcliffe, and Y. Zhan. Integrated proof transformation services. In C. Benz Müller and W. Windsteiger, editors, *IJCAR WS 7*, 2004.