

Partial Regularization of First-Order Resolution Proofs

Jan Gorzny^{a,*}, and Bruno Woltzenlogel Paleo^b

^a *School of Computer Science, University of Waterloo, 200 University Ave. W., Waterloo, ON N2L 3G1, Canada
E-mail: jgorzny@uwaterloo.ca*

^b *College of Engineering and Computer Science, Australian National University, Canberra ACT 0200, Australia
Vienna University of Technology, Karlsplatz 13, 1040, Vienna, Austria
E-mail: bruno@logic.at*

Abstract. This paper describes the generalization of the proof compression algorithm `RecyclePivotsWithIntersection` from propositional to first-order logic. The generalized algorithm performs partial regularization of resolution proofs containing resolution and factoring inferences with *unification*, as generated by many automated theorem provers. An empirical evaluation of the generalized algorithm and its combinations with `GreedyLinearFirstOrderLowerUnits` is also presented.

Keywords: proof compression, first-order logic, resolution, unification

1. Introduction

First-order automated theorem provers, commonly based on resolution and superposition calculi, have recently achieved a high degree of maturity. Proof production is a key feature that has been gaining importance, since proofs are crucial for applications that require certification of a prover's answers or information extractable from proofs (e.g. unsat cores, interpolants, instances of quantified variables). Nevertheless, proof production is non-trivial [12], and the best, most efficient provers do not necessarily generate the best, least redundant proofs.

For proofs using propositional resolution generated by SAT- and SMT-solvers, there is a wide variety of proof compression techniques. Algebraic properties of the resolution operation that might be useful for compression were investigated in [6]. Compression algorithms based on rearranging and sharing chains of resolution inferences have been developed in [2] and [13]. Cotton [5] proposed an algorithm that compresses a

refutation by repeatedly splitting it into a proof of a heuristically chosen literal ℓ and a proof of $\bar{\ell}$, and then resolving them to form a new refutation. The `Reduce&Reconstruct` algorithm [11] searches for locally redundant subproofs that can be rewritten into subproofs of stronger clauses and with fewer resolution steps. A linear time proof compression algorithm based on partial regularization was proposed in [3] and improved in [7].

In contrast, there has been much less work on simplifying first-order proofs. For tree-like sequent calculus proofs, algorithms based on cut-introduction [9, 10] have been proposed. However, converting a DAG-like resolution or superposition proof, as usually generated by current provers, into a tree-like sequent calculus proof may increase the size of the proof. For arbitrary proofs in the TPTP [14] format (including DAG-like first-order resolution proofs), there is a simple algorithm [16] that looks for terms that occur often in any TSTP [14] proof and introduces abbreviations for these terms.

The work reported in this paper is part of a new trend that aims at lifting successful propositional proof compression algorithms to first-order logic. Our first target was the propositional

*Supported by the Google Summer of Code 2014 program. E-mail: jgorzny@uwaterloo.ca.

LowerUnits (LU) algorithm, which delays resolution steps with unit clauses, resulting in the GreedyLinearFirstOrderLowerUnits (GFOLU) algorithm [8]. Here we continue this line of research by lifting the RecyclePivotsWithIntersection (RPI) algorithm [7], which is an improvement of the RecyclePivots (RP) algorithm [3], providing better compression on proofs where nodes have several children.

Section 2 introduces the first-order resolution calculus and the notations used in this paper. Section 4 discusses the challenges that arise in the first-order case (mainly due to unification), which are not present the propositional case. Section 5 describes an algorithm that overcomes these challenges. Section 6 presents experimental results obtained by applying this algorithm, and its combinations with GFOLU, on hundreds of proofs generated with the SPASS theorem prover. Section 7 concludes the paper.

2. The Resolution Calculus

We assume that there are infinitely many variable symbols (e.g. X, Y, Z, X_1, X_2, \dots), constant symbols (e.g. a, b, c, a_1, a_2, \dots), function symbols of every arity (e.g. f, g, f_1, f_2, \dots) and predicate symbols of every arity (e.g. p, q, p_1, p_2, \dots). A *term* is any variable, constant or the application of an n -ary function symbol to n terms. An *atomic formula (atom)* is the application of an n -ary predicate symbol to n terms. A *literal* is an atom or the negation of an atom. The *complement* of a literal ℓ is denoted $\bar{\ell}$ (i.e. for any atom p , $\bar{p} = \neg p$ and $\overline{\neg p} = p$). The set of all literals is denoted \mathcal{L} . A *clause* is a multiset of literals. \perp denotes the *empty clause*. A *unit clause* is a clause with a single literal. Sequent notation is used for clauses (i.e. $p_1, \dots, p_n \vdash q_1, \dots, q_m$ denotes the clause $\{\neg p_1, \dots, \neg p_n, q_1, \dots, q_m\}$). $\text{FV}(t)$ (resp. $\text{FV}(\ell)$, $\text{FV}(\Gamma)$) denotes the set of variables in the term t (resp. in the literal ℓ and in the clause Γ). A *substitution* $\{X_1 \setminus t_1, X_2 \setminus t_2, \dots\}$ is a mapping from variables $\{X_1, X_2, \dots\}$ to, respectively, terms $\{t_1, t_2, \dots\}$. The application of a substitution σ to a term t , a literal ℓ or a clause Γ results in, respectively, the term $t\sigma$, the literal $\ell\sigma$ or the clause $\Gamma\sigma$, obtained from t , ℓ and Γ by replacing all occurrences of the variables in σ by the corresponding terms in σ . The set of all substitutions is denoted \mathcal{S} . A *unifier* of a set of literals is a substitution that makes all literals in the set equal. A *resolution proof* is a directed acyclic graph of clauses where the edges correspond to the inference rules of resolu-

tion and contraction (as explained in detail in Definition 2.1). A *resolution refutation* is a resolution proof with root \perp .

Definition 2.1 (First-Order Resolution Proof).

A directed acyclic graph $\langle V, E, \Gamma \rangle$, where V is a set of nodes and E is a set of edges labeled by literals and substitutions (i.e. $E \subset V \times 2^{\mathcal{L}} \times \mathcal{S} \times V$ and $v_1 \xrightarrow[\sigma]{\ell} v_2$ denotes an edge from node v_1 to node v_2 labeled by the literal ℓ and the substitution σ), is a proof of a clause Γ iff it is inductively constructible according to the following cases:

- **Axiom:** If Γ is a clause, $\hat{\Gamma}$ denotes some proof $\langle \{v\}, \emptyset, \Gamma \rangle$, where v is a new (axiom) node.
- **Resolution:** If ψ_L is a proof $\langle V_L, E_L, \Gamma_L \rangle$ with $\ell_L \in \Gamma_L$ and ψ_R is a proof $\langle V_R, E_R, \Gamma_R \rangle$ with $\ell_R \in \Gamma_R$, and σ_L and σ_R are substitutions such that $\ell_L\sigma_L = \bar{\ell}_R\sigma_R$ and $\text{FV}((\Gamma_L \setminus \{\ell_L\})\sigma_L) \cap \text{FV}((\Gamma_R \setminus \{\ell_R\})\sigma_R) = \emptyset$, then $\psi_L \odot_{\ell_L\sigma_L}^{\sigma_L\sigma_R} \psi_R$ denotes a proof $\langle V, E, \Gamma \rangle$ s.t.

$$V = V_L \cup V_R \cup \{v\}$$

$$E = E_L \cup E_R \cup \left\{ \rho(\psi_L) \xrightarrow[\sigma_L]{\ell_L} v, \rho(\psi_R) \xrightarrow[\sigma_R]{\ell_R} v \right\}$$

$$\Gamma = (\Gamma_L \setminus \{\ell_L\})\sigma_L \cup (\Gamma_R \setminus \{\ell_R\})\sigma_R$$

where v is a new (resolution) node and $\rho(\varphi)$ denotes the root node of φ . The *resolved atom* ℓ is such that $\ell = \ell_L\sigma_L = \bar{\ell}_R\sigma_R$ or $\ell = \bar{\ell}_L\sigma_L = \ell_R\sigma_R$.

- **Contraction:** If ψ' is a proof $\langle V', E', \Gamma' \rangle$ and σ is a unifier of $\{\ell_1, \dots, \ell_n\}$ with $\{\ell_1, \dots, \ell_n\} \subseteq \Gamma'$, then $\lfloor \psi' \rfloor_{\{\ell_1, \dots, \ell_n\}}^\sigma$ denotes a proof $\langle V, E, \Gamma \rangle$ s.t.

$$V = V' \cup \{v\}$$

$$E = E' \cup \{ \rho(\psi') \xrightarrow[\sigma]{\{\ell_1, \dots, \ell_n\}} v \}$$

$$\Gamma = (\Gamma' \setminus \{\ell_1, \dots, \ell_n\})\sigma \cup \{\ell\}$$

where v is a new (contraction) node, $\ell = \ell_k\sigma$ (for any $k \in \{1, \dots, n\}$) and $\rho(\varphi)$ denotes the root node of φ . \square

3. The Propositional Algorithm

RPI (formally defined in Appendix A) removes *irregularities*, which are resolution inferences with a node η when the resolved literal (a.k.a. *pivot*) occurs

as the pivot of another inference located below in the path from η to the root of the proof. In the worst case, regular resolution proofs can be exponentially bigger than irregular ones, but RPI takes care of regularizing the proof only partially, removing inferences only when this does not enlarge the proof.

RPI traverses the proof twice. On the first traversal (bottom-up), it stores for each node a set of *safe literals* that are resolved in all paths below it in the proof or that occur in the root clause of the proof. If one of the node's resolved literals belongs to the set of safe literals, then it is possible to *regularize* the node by replacing it by the parent containing the safe literal. To do this replacement efficiently, the replacement is postponed by marking the other parent as a `deletedNode`. Then, on a single second traversal (top-down), regularization is performed: any node that has a parent node marked as a `deletedNode` is replaced by its other parent.

The RPI and the RP algorithms differ from each other mainly in the computation of the safe literals of a node that has many children. While the former returns the intersection as shown in Algorithm 6, the latter returns the empty set. Moreover, while in RPI the safe literals of the root node contain all the literals of the root clause, in RP the root node is always assigned an empty set of literals.

4. First-Order Challenges

In this section, we describe challenges that have to be overcome in order to successfully adapt RPI to the first-order case. The first example illustrates the need to take unification into account. The other two examples discuss complex issues that can arise when unification is taken into account in a naive way.

Example 4.1. Consider the following proof ψ . When computed as in the propositional case, the safe literals for η_3 are $\{\vdash q(c), p(a, X)\}$.

$$\frac{\eta_1: \vdash p(W, X) \quad \eta_2: p(W, X) \vdash q(c) \quad \eta_4: q(c) \vdash p(a, X)}{\eta_3: \vdash q(c) \quad \eta_5: \vdash p(a, X)} \quad \eta_6: p(Y, b) \vdash \quad \psi: \perp$$

As neither of η_3 's pivots is syntactically equal to a safe literal, the propositional RPI algorithm would not change ψ . However, η_3 's left pivot $p(W, X) \in \eta_1$ is unifiable with the safe literal $p(a, X)$. Regularizing η_3 , by deleting the edge between η_2 and η_3 and replacing

η_3 by η_1 , leads to further deletion of η_4 (because it is not resolvable with η_1) and finally to the much shorter proof below.

$$\frac{\eta_1: \vdash p(W, X) \quad \eta_6: p(Y, b) \vdash}{\psi': \perp}$$

Unlike in the propositional case, where a pivot must be syntactically equal to a safe literal for regularization to be possible, the example above suggests that, in the first-order case, it might suffice that a pivot be unifiable with a safe literal. However, there are cases, as shown in the example below, where mere unifiability is not enough and greater care is needed.

Example 4.2. Again, the safe literals for η_3 , when computed as in the propositional case, are $\{\vdash q(c), p(a, X)\}$, and as the pivot $p(a, c)$ is unifiable with the safe literal $p(a, X)$, η_3 appears to be a candidate for regularization.

$$\frac{\eta_1: \vdash p(a, c) \quad \eta_2: p(a, c) \vdash q(c) \quad \eta_4: q(c) \vdash p(a, X)}{\eta_3: \vdash q(c) \quad \eta_5: \vdash p(a, X)} \quad \eta_6: p(Y, b) \vdash \quad \psi: \perp$$

However, if we attempt to regularize the proof, the same series of actions as in Example 4.1 would require resolution between η_1 and η_6 , which is not possible.

One way to prevent the problem depicted above would be to require the pivot to be not only unifiable but in fact more general than a safe literal. A weaker (and better) requirement is possible, however, as defined below.

Definition 4.1. Let η be a node with pivot ℓ' unifiable with safe literal ℓ which is resolved against literals ℓ_1, \dots, ℓ_n in a proof ψ . η is said to satisfy the *pre-regularization unifiability property* in ψ if ℓ_1, \dots, ℓ_n , and $\bar{\ell}'$ are unifiable.

One way to ensure this property is met is to slightly modify the notion of safe literals, by applying the unifier of the resolution step to the each pivot before adding it to the safe literals (cf. algorithm 3, lines 8 and 10). In the case of Example 4.2, this would result in η_3 having the safe literals $\{\vdash q(c), p(a, b)\}$, where clearly the pivot $p(a, c)$ in η_1 is not safe.

Example 4.3. Satisfying the pre-regularization unifiability property is not sufficient. Consider the proof ψ in Figure 1. After collecting the safe literals, η_3 's safe literals are $\{q(T, V), p(c, d) \vdash q(f(a, e), c)\}$. η_3 's pivot $q(f(a, V), U)$ is unifiable to (and even more general

$$\begin{array}{c}
\frac{\eta_1: p(U, V) \vdash q(f(a, V), U) \quad \eta_2: q(f(a, X), Y), q(T, X) \vdash q(f(a, Z), Y)}{\eta_3: p(U, V), q(T, V) \vdash q(f(a, Z), U)} \quad \eta_4: \vdash q(R, S) \\
\frac{\eta_6: \vdash p(c, d) \quad \eta_5: p(U, V) \vdash q(f(a, Z), U)}{\eta_7: \vdash q(f(a, Z), c)} \\
\frac{\eta_8: q(f(a, e), c) \vdash \quad \eta_7: \vdash q(f(a, Z), c)}{\psi: \perp}
\end{array}$$

Fig. 1. An example where pre-regularization unifiability is not sufficient.

than) the safe literal $q(f(a, e), c)$. Attempting to regularize η_3 would lead to the removal of η_2 , the replacement of η_3 by η_1 and the removal of η_4 (because η_1 does not contain the pivot required by η_5), with η_5 also being replaced by η_1 . Then resolution between η_1 and η_6 results in η'_7 , which cannot be resolved with η_8 , as shown below.

$$\frac{\eta_6: \vdash p(c, d) \quad \eta_1: p(U, V) \vdash q(f(a, V), U)}{\eta'_7: \vdash q(f(a, d), c)} \quad \eta_8: q(f(a, e), c) \vdash$$

$\psi': ??$

η_1 's literal $q(f(a, V), U)$, which would be resolved with η_8 's literal, was changed to $q(f(a, d), c)$ due to the resolution between η_1 and η_6 .

Thus we additionally require the following property be satisfied.

Definition 4.2. Let η be a node with safe literals ϕ that is marked for regularization with parents η_1 and η_2 , where η_2 is marked as a `deletedNode` in a proof ψ . η is said to satisfy the *regularization unifiability property* in ψ if there exists a substitution σ such that $\eta_1\sigma \subseteq \phi$.

This property ensures that the remainder of the proof does not expect a variable in η_1 to be unified to different values simultaneously. This property is not necessary in the propositional case, as the replacement node would not change lower in the proof.

5. First-Order RecyclePivotsWithIntersection

This section presents `FirstOrderRecyclePivotsWithIntersection` (FORPI), Algorithm 1, a first-order generalization of RPI. FORPI traverses the proof in a bottom-up manner, storing for every node a set of safe literals. The set of safe literals for a node ψ is computed from the set of safe literals of its children (cf. Algorithm 3), similarly to the propositional case, but additionally applying unifiers to the resolved pivots (cf. Example 4.2). If one of the node's resolved literals can be unified to a literal in the set of

input : A first-order proof ψ
output: A possibly less-irregular first-order proof ψ'

```

1  $\psi' \leftarrow \psi$ ;
2 traverse  $\psi'$  bottom-up and foreach node  $\eta$  in  $\psi'$  do
3   if  $\eta$  is a resolvent node then
4     setSafeLiterals( $\eta$ ) ;
5     regularizeIfPossible( $\eta$ )
6  $\psi' \leftarrow \text{fix}(\psi')$  ;
7 return  $\psi'$ ;

```

Algorithm 1: FORPI

input : A node $\psi = \psi_L \odot_{\ell_L \ell_R}^{\sigma_L \sigma_R} \psi_R$
output: nothing (but the proof containing ψ may be changed)

```

1 if  $\exists \sigma$  and  $\ell \in \psi.\text{safeLiterals}$  such that  $\ell\sigma = \ell_R$  or  $\ell = \ell_R\sigma$  then
2   if  $\exists \sigma'$  such that  $\psi_R\sigma' \subseteq \psi.\text{safeLiterals}$  then
3     mark  $\psi_L$  as deletedNode ;
4     mark  $\psi$  as regularized
5 else if  $\exists \sigma$  and  $\ell \in \psi.\text{safeLiterals}$  such that  $\ell\sigma = \ell_L$  or  $\ell = \ell_L\sigma$  then
6   if  $\exists \sigma'$  such that  $\psi_L\sigma' \subseteq \psi.\text{safeLiterals}$  then
7     mark  $\psi_R$  as deletedNode ;
8     mark  $\psi$  as regularized

```

Algorithm 2: FRegularizeIfPossible

safe literals, then it may be possible to regularize the node by replacing it by one of its parents.

In the first-order case, we additionally check for the regularization property (cf. lines 2 and 6 of Algorithm 2). Similarly to RPI, instead of replacing the irregular node by one of its parents immediately, its other parent is marked as a `deletedNode`, as shown in Algorithm 2. As in the propositional case, fixing of the proof is postponed to another (single) traversal, as regularization proceeds top-down and only nodes below a regularized node may require fixing. During fixing, the irregular node is actually replaced by the parent that is not marked as `deletedNode`. During proof fixing, factoring inferences can be applied, in order to compress the proof further.

```

input : A first-order resolution node  $\psi$ 
output: nothing (but the node  $\psi$  gets a set of safe literals)

1 if  $\psi$  is a root node with no children then
2    $\psi.\text{safeLiterals} \leftarrow \psi.\text{clause}$ 
3 else
4   foreach  $\psi' \in \psi.\text{children}$  do
5     if  $\psi'$  is marked as regularized then
6        $\text{safeLiteralsFrom}(\psi') \leftarrow \psi'.\text{safeLiterals}$ ;
7     else if  $\psi' = \psi \odot_{\ell_L \ell_R}^{\sigma_L \sigma_R} \psi_R$  for some  $\psi_R$  then
8        $\text{safeLiteralsFrom}(\psi') \leftarrow \psi'.\text{safeLiterals} \cup \{ \ell_R \sigma_R \}$ 
9     else if  $\psi' = \psi_L \odot_{\ell_L \ell_R}^{\sigma_L \sigma_R} \psi$  for some  $\psi_L$  then
10       $\text{safeLiteralsFrom}(\psi') \leftarrow \psi'.\text{safeLiterals} \cup \{ \ell_L \sigma_L \}$ 
11    $\psi.\text{safeLiterals} \leftarrow \bigcap_{\psi' \in \psi.\text{children}} \text{safeLiteralsFrom}(\psi')$ 

```

Algorithm 3: FOsetSafeLiterals

6. Experiments

A prototype version of FORPI has been implemented in the functional programming language Scala as part of the Skeptik library. This library includes an implementation of GFOLU [8]. In order to evaluate the algorithm's effectiveness, FORPI was tested on two data sets: proofs generated by a real theorem prover and randomly-generated resolution proofs. The proofs are included in the source code repository, available at <https://github.com/jgorzny/Skeptik>. Note that by implementing the algorithms in this library, we are able to guarantee the correctness of the compressed proofs, as Skeptik is also a proof-checker.

First, FORPI was evaluated on the same proofs used to evaluate GFOLU. This data was generated by executing the SPASS (<http://www.spass-prover.org/>) theorem prover on 2280 real first-order problems without equality of the TPTP Problem Library (among them, 1032 problems are known to be unsatisfiable). In order to generate pure resolution proofs, the advanced inference rules of SPASS were disabled. The proofs were originally generated on the Euler Cluster at the University of Victoria with a time limit of 300 seconds per problem. Under these conditions, SPASS generated 308 proofs. The proofs generated by SPASS were small: proof lengths varied from 3 to 49, and the number of resolutions in a proof ranged from 1 to 32.

In order to test FORPI's effectiveness on larger proofs, randomly generated first-order resolution proofs were also used. Proofs were generated by the following procedure: start with a root node whose conclusion is \perp , and make two premises η_1 and η_2 using a randomly generated literal such that the desired con-

clusion is the result of resolving η_1 and η_2 . For each node η_i , determine the inference rule used to make its conclusion: with probability $p = 0.9$, η_i is the result of a resolution, otherwise it is the result of a contraction. Literals are generated by uniformly choosing a predicate from the set of $\{1, \dots, k+1\}$ where k is the number of predicates generated so far; the choice $k+1$ generates a new predicate with a new random arity (at most four). For each argument, it is a constant with probability $p = 0.7$ and a function otherwise; functions are generated similarly to predicates. If a node η should be the result of a resolution, then with probability $p = 0.2$ we generate a left parent η_ℓ and a right parent η_r for η (i.e. $\eta = \eta_\ell \odot \eta_r$) having a common parent η_c (i.e. $\eta_\ell = (\eta_\ell)_\ell \odot \eta_c$ and $\eta_r = \eta_c \odot (\eta_r)_r$, for some newly generated nodes $(\eta_\ell)_\ell$ and $(\eta_r)_r$). The common parent ensures that also non-tree-like DAG proofs are generated. This procedure is recursively applied to the generated parent nodes. Each parent of a resolution has each of its constants and functions replaced by a fresh variable with probability $p = 0.7$. At each recursive call, the additional minimum height required for the remainder of the branch is decreased by one with probability $p = 0.5$. Thus if each branch always decreases the additional required height, the proof has height equal to the initial minimum value. The process stops when every branch is required to add a subproof of height zero, or after a timeout. The end of each branch is taken as an axiom. The minimum height was set to 7 (which is the minimum number of nodes in an irregular proof plus one) and the timeout was set to 300 seconds (the same timeout allowed for SPASS).

In order to match the maximum number of potential proofs from the TPTP Problem Library, 2280 randomly generated proofs were used as the second data set. The randomly generated proofs were much larger than those of the first data set: proof lengths varied from 95 to 700, while the number of resolutions in a proof ranged from 48 to 368.

For consistency, the same system and metrics were used. Proof compression and proof generation was performed on a laptop (2.8GHz Intel Core i7 processor with 4GB of RAM (1333MHz DDR3) available to the Java Virtual Machine). For each proof ψ , we measured the time needed to compress the proof ($t(\psi)$) and the compression ratio in terms of resolutions $((|\psi| - |\alpha(\psi)|)/|\psi|)$ where $|\psi|$ is the number of resolutions in the proof, and $\alpha(\psi)$ is the result of applying a compression algorithm or some composition of FORPI and GFOLU. Note that we consider only

the number of resolutions in order to compare the results of these algorithms to their propositional variants (where contractions are implicit). Moreover, contractions could be made implicit within resolution inferences even in the first-order case and we use explicit contractions only for technical convenience.

Table 1 summarizes the results of `FORPI` and its combinations with `GFOLU`. The first set of columns describes the percentage of proofs that were compressed by each compression algorithm. The algorithm ‘Best’ runs both of combinations of `GFOLU` and `FORPI` and returns the shortest proof output by either of them. The total number of proofs is $308 + 2280 = 2588$ and the total number of resolution nodes is $2,249 + 393,883 = 396,132$. The percentages in the last three columns are computed by $(\sum_{\psi \in \Psi} |\psi| - \sum_{\psi \in \Psi} |\alpha(\psi)|) / (\sum_{\psi \in \Psi} |\psi|)$ for each data set Ψ (TPTP, Random, or Both). The use of `FORPI` alongside `GFOLU` allows at least an additional 5% of proofs to be compressed. Furthermore, the use of both algorithms removes more than twice as many nodes than any single algorithm.

Table 2 compares the results of `FORPI` and its combinations with `GFOLU` with their propositional variants as evaluated in [4]. The first column describes the mean compression ratio for each algorithm including proofs that were not compressed by the algorithm, while the second column calculates the mean compression ratio considering only compressed proofs. It is unsurprising that the first column is lower than the propositional mean for each algorithm: there are stricter requirements to apply these algorithms to first-order proofs. In particular, additional properties must be satisfied before a unit can be lowered, or before a pivot can be recycled. On the other hand, when first-order proofs are compressed, the levels of compression are on par with or better than their propositional counterparts.

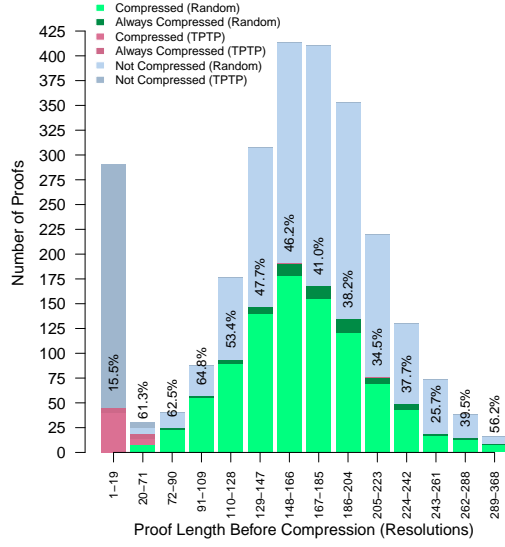
Figure 2 (a) shows the number of proofs (compressed and uncompressed) per grouping based on number of resolutions in the proof. The red (resp. dark grey) data shows the number of compressed (resp. uncompressed) proofs for the TPTP data set, while the green (resp. light grey) data shows the number of compressed (resp. uncompressed) proofs for the random proofs. The number of proofs in each group is the sum of the heights of each coloured bar in to that group. The overall percentage of proofs compressed in a group is indicated on each bar. Dark colors indicate the number of proofs compressed by `FORPI`, `GFOLU`, and both compositions of these algorithms; light colors indicate cases where `FORPI` succeeded, but at least one of

`GFOLU` or a combination of these algorithms achieved zero compression. Given the size of the TPTP proofs, it is unsurprising that few are compressed: small proofs are a priori less likely to contain irregularities. On the other hand, at least 25% of randomly generated proofs in each size group could be compressed.

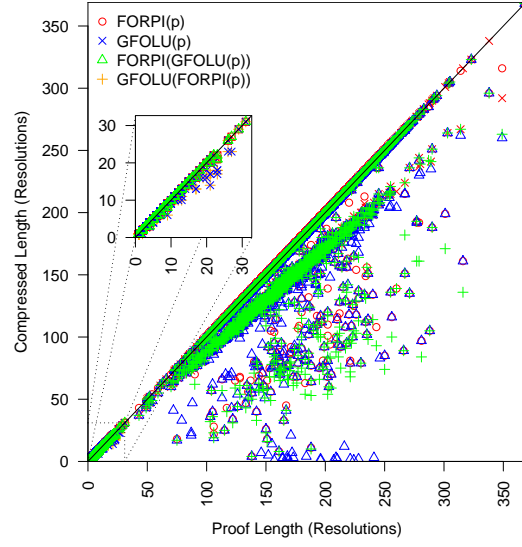
Figure 2 (b) is a scatter plot comparing the number of resolutions of the input proof against the number of resolutions in the compressed proof for each algorithm. The results on the TPTP data are magnified in the sub-plot. For the randomly generated proofs (points outside of the sub-plot), it is often the case that the compressed proof is significantly shorter than the input proof. Interestingly, `GFOLU` appears to reduce the number of resolutions by a linear factor in many cases. This is likely due to a linear growth in the number of non-interacting irregularities (i.e. irregularities for which the lowered units share no common literals with any other sub-proofs), which leads to a linear number of nodes removed.

Figure 2 (c) is a scatter plot comparing the size of compression obtained by applying `FORPI` before `GFOLU` versus `GFOLU` before `FORPI`. Data obtained from the TPTP data set is marked in red; the remaining points are obtained from randomly generated proofs. Points that lie on the diagonal line have the same size after each combination. There are 165 points beneath the line and 258 points above the line. Therefore, as in the propositional case [7], it is not a priori clear which combination is better. Nevertheless, the distinctly greater number of points above the line suggests that it is more often the case that `FORPI` should be applied after `GFOLU`. Not only this combination is not only more likely to maximize the likelihood of compression, but the achieved compression also tends to be larger.

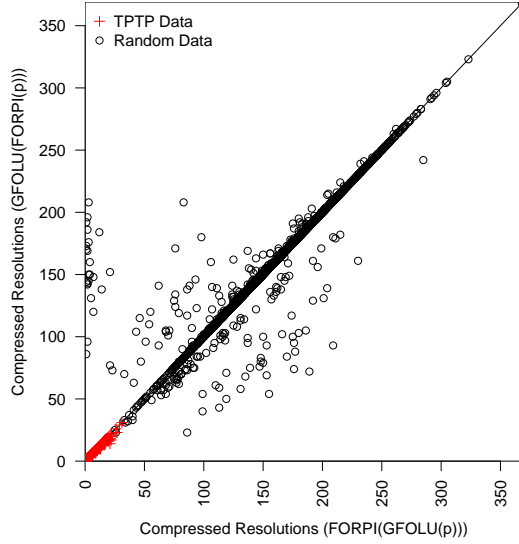
Figure 2 (d) shows a plot comparing the difference between the cumulative number of resolutions of the first x input proofs and the cumulative number of resolutions in the first x proofs after compression (i.e. the cumulative number of *removed* resolutions). The TPTP data is displayed in the sub-plot; note that the lines for everything except `FORPI` largely overlap (since the values are almost identical; cf. Table 1). Observe that the use of both algorithms is always better than using a single algorithm. The data also shows that using `FORPI` after `GFOLU` is normally the preferred order of composition, as it typically results in a greater number of nodes removed than the other combination. An even better approach is to try both com-



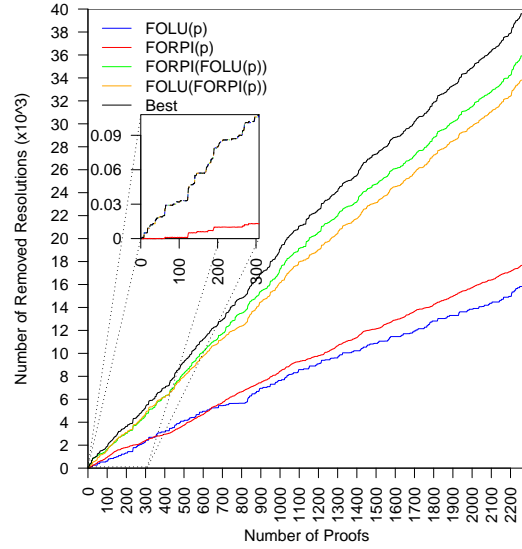
(a) Number of (non)-compressed proofs



(b) Compressed length against input length



(c) FORPI (GFOLU (p)) vs. GFOLU (FORPI (p))



(d) Cumulative proof compression

Fig. 2. GFOLU & FORPI Combination Results

Algorithm	# of Proofs Compressed			# of Removed Nodes		
	TPTP	Random	Both	TPTP	Random	Both
GFOLU(p)	55 (17.9%)	815 (35.7%)	870 (33.6%)	107 (4.8%)	17,730 (4.5%)	17,837 (4.5%)
FORPI(p)	11 (3.6%)	252 (11.1%)	263 (10.2%)	13 (0.6%)	15,913 (4.0%)	15,926 (4.0%)
GFOLU(FORPI(p))	55 (17.9%)	993 (43.6%)	1048 (40.5%)	108 (4.8%)	33,956 (9.1%)	34,064 (9.1%)
FORPI(GFOLU(p))	11 (3.6%)	993 (43.6%)	1004 (38.8%)	108 (4.8%)	36,070 (9.1%)	36,178 (9.1%)
Best	56 (18.2%)	993 (43.6%)	1049 (40.5%)	108 (4.8%)	39,742 (10.1%)	39,850 (10.1%)

Table 1

Number of proofs compressed and number of overall nodes removed

Algorithm	First-Order Compression		Algorithm	Propositional Compression [4]
	All	Compressed Only		
GFOLU(p)	3.4%	33.1%	LU(p)	7.5%
FORPI(p)	4.5%	13.4%	RPI(p)	17.8%
GFOLU(FORPI(p))	7.6%	19.7%	(LU(RPI(p)))	21.7%
FORPI(GFOLU(p))	8.1%	21.0%	(RPI(LU(p)))	22.0%
Best	9.2%	22.8%	Best	22.0%

Table 2

Mean compression results

binations and choose the best result (as shown in the ‘Best’ curve).

SPASS required approximately 40 minutes (running on a cluster and including proof generation time for each problem) to generate all the 308 TPTP proofs. The total time to apply both FORPI and GFOLU on all these proofs was just over 8 seconds on a simple laptop computer. The random proofs were generated in 70 minutes, and took approximately 453 seconds (or 7.5 minutes) to compress, both measured on the same computer. All times include parsing time. These compression algorithms continue to be very fast, and may simplify the proof considerably for a relatively small cost in time.

7. Conclusions and Future Work

The main contribution of this paper is the generalization of the propositional proof compression algorithm RPI to the first-order case. As indicated in Section 4, the generalization is challenging, because unification changes the pivots and, consequently, must be taken into account when collecting safe literals and marking nodes for deletion.

Every computational experiment evaluates not only the algorithm but also the data on which it is executed. Although the experimental results are not as promising as expected, this is due to the fact that the 308 proofs currently available are too short to con-

tain a significant amount of irregularities. This is a valuable piece of information, allowing us to conclude that it is not worth applying FORPI to pure resolution proofs which current state-of-the-art first-order theorem provers seem capable of producing. Nevertheless, based on our positive results for RPI on much longer proofs generated by SAT and SMT solvers [7], FORPI remains a promising option to be revisited in the future, when the performance of first-order theorem provers catch up with advances in SAT and SMT and taller first-order benchmark proofs become available.

References

- [1] *Logic for Programming, Artificial Intelligence, and Reasoning 16th International Conference, Dakar, Senegal, Revised Selected Papers*, LNCS. Springer, 2010.
- [2] H. Amjad. Compressing propositional refutations. *Electronic Notes in Theoretical Computer Science*, 185:3–15, 2007.
- [3] O. Bar-Ilan, O. Fuhrmann, S. Hoory, O. Shacham, and O. Strichman. Linear-time reductions of resolution proofs. In *Haifa Verification Conference*, LNCS, pages 114–128. Springer, 2008.
- [4] J. Boudou and B. Woltzenlogel Paleo. Compression of propositional resolution proofs by lowering subproofs. In *Automated Reasoning with Analytic Tableaux and Related Methods - 22th International Conference*, LNCS, pages 59–73. Springer, 2013.
- [5] S. Cotton. Two techniques for minimizing resolution proofs. In Ofer Strichman and Stefan Szeider, editors, *SAT 2010*, LNCS, pages 306–312. Springer, 2010.

- [6] P. Fontaine, S. Merz, and B. Woltzenlogel Paleo. Exploring and exploiting algebraic and graphical properties of resolution. In *8th International Workshop on SMT*, 2010.
- [7] P. Fontaine, S. Merz, and B. Woltzenlogel Paleo. Compression of propositional resolution proofs via partial regularization. In *Automated Deduction - CADE-23 - 23rd International Conference on Automated Deduction*, Wroclaw, Poland, July 31 - August 5, 2011. *Proceedings*, LNCS, pages 237–251. Springer, 2011.
- [8] J. Gorzny and B. Woltzenlogel Paleo. Towards the compression of first-order resolution proofs by lowering unit clauses. In *Automated Deduction - CADE-25 - 25th International Conference on Automated Deduction*, Berlin, Germany, August 1-7, 2015. *Proceedings*, 2015.
- [9] S. Hetzl, A. Leitsch, G. Reis, and D. Weller. Algorithmic introduction of quantified cuts. *Theoretical Computer Science*, 549:1–16, 2014.
- [10] B. Woltzenlogel Paleo. Atomic cut introduction by resolution: Proof structuring and compression. In *LPAR-16* [1], pages 463–480.
- [11] S. F. Rollini, R. Bruttomesso, and N. Sharygina. An efficient and flexible approach to resolution proof reduction. In *Hardware and Software: Verification and Testing*, LNCS, pages 182–196. Springer, 2011.
- [12] S. Schulz and G. Sutcliffe. Proof generation for saturating first-order theorem provers. In D. Delahaye and B. Woltzenlogel Paleo, editors, *All about Proofs, Proofs for All*, volume 55 of *Mathematical Logic and Foundations*. College Publications, London, UK, 2015.
- [13] C. Sinz. Compressing propositional proofs by common subproof extraction. In R. Moreno-Díaz, F. Pichler, and A. Quesada-Arencibia, editors, *EUROCAST*, LNCS, pages 547–555. Springer, 2007.
- [14] G. Sutcliffe. The TPTP Problem Library and Associated Infrastructure: The FOF and CNF Parts, v3.5.0. *Journal of Automated Reasoning*, 43(4):337–362, 2009.
- [15] G. S. Tseitin. On the complexity of derivation in propositional calculus. In J. Siekmann and G. Wrightson, editors, *Automation of Reasoning: Classical Papers in Computational Logic 1967-1970*. Springer-Verlag, 1983.
- [16] J. Vyskocil, D. Stanovský, and J. Urban. Automated proof compression by invention of new definitions. In *LPAR* [1], pages 447–462.

Appendix A. Algorithm

RecyclePivotsWithIntersection

Note: for the reviewers' convenience, this appendix summarizes [7].

RecyclePivotsWithIntersection (RPI) [7] aims at compressing irregular proofs. It can be seen as a simple but significant modification of the RP algorithm described in [3], from which it derives its name. Although in the worst case full regularization can increase the proof length exponentially [15], these algo-

rithms show that many irregular proofs can have their length decreased if a careful partial regularization is performed.

Consider an irregular proof of the form $\psi[\eta \odot_p \psi'[\eta' \odot_p \eta'']]$ and assume, without loss of generality, that $p \in \eta$ and $p \in \eta'$. Then, if $\eta' \odot_p \eta''$ is replaced by η'' within the proof-context $\psi'[\]$, the clause $\eta \odot_p \psi'[\eta'']$ subsumes the clause $\eta \odot_p \psi'[\eta' \odot_p \eta'']$, because even though the literal $\neg p$ of η'' is propagated down, it gets resolved against the literal p of η later on below in the proof. More precisely, even though it might be the case that $\neg p \in \psi'[\eta'']$ while $\neg p \notin \psi'[\eta' \odot_p \eta'']$, it is necessarily the case that $\neg p \notin \eta \odot_p \psi'[\eta' \odot_p \eta'']$ and $\neg p \notin \eta \odot_p \psi'[\eta'']$.

Although the remarks above suggest that it is safe to replace $\eta' \odot_p \eta''$ by η'' within the proof-context $\psi'[\]$, this is not always the case. If a node in $\psi'[\]$ has a child in $\psi[\]$, then the literal $\neg p$ might be propagated down to the root of the proof, and hence, the clause $\psi[\eta \odot_p \psi'[\eta'']]$ might not subsume the clause $\psi[\eta \odot_p \psi'[\eta' \odot_p \eta'']]$. Therefore, it is only safe to do the replacement if the literal $\neg p$ gets resolved in all paths from η'' to the root or if it already occurs in the root clause of the original proof $\psi[\eta \odot_p \psi'[\eta' \odot_p \eta'']]$.

These observations lead to the idea of traversing the proof in a bottom-up manner, storing for every node a set of *safe literals* that get resolved in all paths below it in the proof (or that already occurred in the root clause of the original proof). Moreover, if one of the node's resolved literals belongs to the set of safe literals, then it is possible to regularize the node by replacing it by one of its parents (cf. Algorithm 4).

The regularization of a node should replace a node by one of its parents, and more precisely by the parent whose clause contains the resolved literal that is safe. After regularization, all nodes below the regularized node may have to be fixed. However, since the regularization is done with a bottom-up traversal, and only nodes below the regularized node need to

<p>input : A proof ψ output: A possibly less-irregular proof ψ'</p> <pre> 1 $\psi' \leftarrow \psi$; 2 traverse ψ' bottom-up and foreach node η in ψ' do 3 if η is a resolvent node then 4 $\text{setSafeLiterals}(\eta)$; 5 $\text{regularizeIfPossible}(\eta)$ 6 $\psi' \leftarrow \text{fix}(\psi')$; 7 return ψ'; </pre>
--

Algorithm 4: RPI

be fixed, it is again possible to postpone fixing and do it with only a single traversal afterwards. Therefore, instead of replacing the irregular node by one of its parents immediately, its other parent is marked as `deletedNode`, as shown in Algorithm 5. Only later during fixing, the irregular node is actually replaced by its surviving parent (i.e. the parent that is not marked as `deletedNode`).

The set of safe literals of a node η can be computed from the set of safe literals of its children (cf. Algorithm 6). In the case when η has a single child ς , the safe literals of η are simply the safe literals of ς together with the resolved literal p of ς belonging to η (p is safe for η , because whenever p is propagated down the proof through η , p gets resolved in ς). It is important to note, however, that if ς has been marked as regularized, it will eventually be replaced by η , and hence p should not be added to the safe literals of η . In this case, the safe literals of η should be exactly the same as the safe literals of ς . When η has several children, the safe literals of η w.r.t. a child ς_i contain literals that are safe on all paths that go from η through ς_i to the root. For a literal to be safe for all paths from η to the root, it should therefore be in the intersection of the sets of safe literals w.r.t. each child.

The RP and the RPI algorithms differ from each other mainly in the computation of the safe literals of a node that has many children. While RPI returns the intersection as shown in Algorithm 6, RP returns the empty set (cf. Algorithm 7). Additionally, while in RPI the safe literals of the root node contain all the literals of the root clause, in RP the root node is always assigned an empty set of literals. (Of course, this makes a difference only when the proof is not a refutation.) Note that during a traversal of the proof, the lines from 5 to 10 in Algorithm 6 are executed as many times as the number of edges in the proof. Since every node has at most two parents, the number of edges is at most twice the number of nodes. Therefore, during a traversal of a proof with n nodes, lines from 5 to 10 are executed at most $2n$ times, and the algorithm remains linear. In our prototype implementation, the sets of safe literals are instances of Scala's `mutable.HashSet` class. Being mutable, new elements can be added efficiently. And being HashSets, membership checking is done in constant time in the average case, and set in-

tersection (line 12) can be done in $O(k.s)$, where k is the number of sets and s is the size of the smallest set.

```

input : A node  $\eta$ 
output: nothing (but the proof containing  $\eta$  may be changed)

1 if  $\eta$ .rightResolvedLiteral  $\in \eta$ .safeLiterals then
2   mark left parent of  $\eta$  as deletedNode ;
3   mark  $\eta$  as regularized
4 else if  $\eta$ .leftResolvedLiteral  $\in \eta$ .safeLiterals then
5   mark right parent of  $\eta$  as deletedNode ;
6   mark  $\eta$  as regularized

```

Algorithm 5: `regularizeIfPossible`

```

input : A node  $\eta$ 
output: nothing (but the node  $\eta$  gets a set of safe literals)

1 if  $\eta$  is a root node with no children then
2    $\eta$ .safeLiterals  $\leftarrow \eta$ .clause
3 else
4   foreach  $\eta' \in \eta$ .children do
5     if  $\eta'$  is marked as regularized then
6       safeLiteralsFrom( $\eta'$ )  $\leftarrow \eta'$ .safeLiterals ;
7     else if  $\eta$  is left parent of  $\eta'$  then
8       safeLiteralsFrom( $\eta'$ )  $\leftarrow \eta'$ .safeLiterals  $\cup \{ \eta'.rightResolvedLiteral \}$  ;
9     else if  $\eta$  is right parent of  $\eta'$  then
10      safeLiteralsFrom( $\eta'$ )  $\leftarrow \eta'$ .safeLiterals  $\cup \{ \eta'.leftResolvedLiteral \}$  ;
11  $\eta$ .safeLiterals  $\leftarrow \bigcap_{\eta' \in \eta.children} \text{safeLiteralsFrom}(\eta')$ 

```

Algorithm 6: `setSafeLiterals`

```

input : A node  $\eta$ 
output: nothing (but the node  $\eta$  gets a set of safe literals)

1 if  $\eta$  is a root node with no children then
2    $\eta$ .safeLiterals  $\leftarrow \emptyset$ 
3 else
4   if  $\eta$  has only one child  $\eta'$  then
5     if  $\eta'$  is marked as regularized then
6        $\eta$ .safeLiterals  $\leftarrow \eta'$ .safeLiterals ;
7     else if  $\eta$  is left parent of  $\eta'$  then
8        $\eta$ .safeLiterals  $\leftarrow \eta'$ .safeLiterals  $\cup \{ \eta'.rightResolvedLiteral \}$  ;
9     else if  $\eta$  is right parent of  $\eta'$  then
10       $\eta$ .safeLiterals  $\leftarrow \eta'$ .safeLiterals  $\cup \{ \eta'.leftResolvedLiteral \}$  ;
11   else
12      $\eta$ .safeLiterals  $\leftarrow \emptyset$ 

```

Algorithm 7: `setSafeLiterals` for RP