

Lemmas: Generation, Selection, Application^{*}

Michael Rawson¹, Christoph Wernhard², Zsolt Zombori³, and
Wolfgang Bibel⁴

¹ TU Wien, Austria michael@rawsons.uk

² University of Potsdam, Germany info@christophwernhard.com

³ Alfréd Rényi Institute of Mathematics, Hungary zombori@renyi.hu

⁴ Technical University Darmstadt, Germany bibel@gmx.net

Abstract. Noting that lemmas are a key feature of mathematics, we engage in an investigation of the role of lemmas in automated theorem proving. The paper describes experiments with a combined system involving learning technology that generates useful lemmas for automated theorem provers, demonstrating improvement for several representative systems and solving a hard problem not solved by any system for twenty years. By focusing on condensed detachment problems we simplify the setting considerably, allowing us to get at the essence of lemmas and their role in proof search.

1 Introduction

Mathematics is built in a carefully structured way, with many disciplines and subdisciplines. These are characterized by concepts, definitions, axioms, theorems, lemmas, and so forth. There is no doubt that this inherent structure of mathematics is part of the discipline's long-lasting success.

Research into Automated Theorem Proving (ATP) to date has taken little notice of the information provided by this structure. Even state-of-the-art ATP systems ingest a conjecture together with pertinent definitions and axioms in a way completely agnostic to their place in the mathematical structure. A comparatively small but nevertheless important part of the structure of mathematics is the identification and application of *lemmas*. It is this aspect which is the focus of the work presented here.

The purpose of lemmas in mathematics is at least threefold. First, and perhaps most importantly, lemmas support the search for proofs of assertions. If some lemma applies to a given problem, a proof may be found more easily. Second, it is often the case that a lemma may be applied more than once. If this happens, its use will shorten the length of the overall proof since the proof of the lemma need only be carried out once, not repeatedly for every application.

^{*} Funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – Project-ID 457292495, by the North-German Supercomputing Alliance (HLRN), by the ERC grant CoG ARTIST 101002685, by the Hungarian National Excellence Grant 2018-1.2.1-NKP-00008 and the Hungarian Artificial Intelligence National Laboratory Program (RRF-2.3.1-21-2022-00004).

Third, the structuring effect of proofs by the use of lemmas is an important feature for human comprehension of proofs. In our work we are motivated primarily by the first two of these three aspects.

These considerations give rise to the crucial question: how can we find useful lemmas for proving a given problem? Here we mean useful in the sense of the two aforementioned aspects: lemmas should be applicable to the problem at hand, preferably many times. In full generality this is a difficult question indeed, which will require much further research. In this first step we restrict the question to a narrow range of problems, known in literature as *condensed detachment* (CD) problems [39]. Our investigation therefore focuses on answers to the question of how ATP performance may be improved for CD problems by the generation and selection of useful lemmas before search begins.

CD problems are of the form “axiom(s) and *Det* imply a goal” where *Det* represents the well-known modus ponens rule, or condensed detachment. In order to support this study experimentally, we have built a combined system for dealing with these problems. It features SGCD [72] as prover and lemma generator along with a learning module based on either an easily-interpreted linear model over novel hand-engineered features, or a graph neural network supporting end-to-end learning directly from lemmas.

We experiment⁵ with a total of 312 CD problems, including all 196 pure CD problems from TPTP [63]. We employ several representative ATP systems, including state-of-the-art systems for both general first-order reasoning and for CD problems. SGCD solves in our combined system with learned lemmas 287 of the 312 problems. These include 28 problems not solved by the leading first-order prover Vampire [27], which solves 263 problems in its *CASC* [62] portfolio mode. Supplemented with our lemmas, Vampire is boosted to 284 solved problems. In combination, SGCD and the boosted Vampire give 293 solved problems. Taking into account the solutions obtained by further provers with our lemmas, we obtain a total of 296 — plus one more discussed next.

An outstanding success of the approach taken here is the automatic proof of the Meredith single axiom theorem LCL073-1, which has persisted in TPTP with rating 1.00 and status Unknown since 1997. The first and only system to succeed was OTTER [37], after intensive massaging by Wos [83]. With SGCD, finding a proof is straightforward. Apparently, the combination of the Connection Method (CM) [5,6,7] and our approach to lemmas is a powerful one.

Section 2 presents condensed detachment and its embedding into the CM by way of so-called *D-terms*, as well as background material on lemmas and machine learning in ATP. Section 3 introduces a method for generating and selecting useful lemmas and presents experimental results with it, leading up to the proof of LCL073-1 in Section 4. We conclude with a detailed summary and outlook for further work in this area in Section 5.

⁵ All experiments are fully reproducible and the artefacts are available at <https://github.com/zsoltzombori/lemma>, commit `df2faaa`. We use CD Tools [72], PIE [70,71] and SWI-Prolog [76] for reasoning tasks, and PyTorch [46] for learning.

2 Background and Related Work

In a very general sense, lemmas in ATP factorize duplication. This may be between different proofs that make use of the same lemma, or within a single proof where a lemma is used multiple times. It may not even be a particular formula that is shared, but a *pattern*, such as a *resonator* [80]. In the presence of machine learning, we may think of even more abstract entities that are factorized: the *principles* by which proofs are written, repeated in different proofs or contexts.

Depending on the proving method, lemmas in ATP play different roles. Provers based on *saturation*, typically resolution/superposition (RS) systems [3], inherently operate by generating lemmas: a resolvent is itself a lemma derived from its parents. Nevertheless, one may ask for more meaningful lemmas than the clauses of the proof. This is addressed with *cut introduction* [77,18,11], which studies methods to obtain complex lemmas from resolution proofs. Such lemmas provide insight about the high-level structure of proofs, extract interesting concepts and support research into the correspondence between natural mathematical notions and possible proof compressions. Other approaches to interesting theorems or lemmas are described for example in [64,51].

Another question concerning lemmas and ATP systems is whether performance can be improved by supplementing the input with lemmas. This is particularly applicable if lemmas were obtained with *different* methods than those of the prover. Otherwise, it may have obtained these by itself. We note here that in some cases systems *cannot* generate certain lemmas because of e.g. ordering restrictions. As we will see, leading ATP systems such as Vampire and E [58] can indeed be improved in this way. *Different methods* does not necessarily mean *different systems*: it is possible to use different configurations of the same system for lemma generation and proving, as well as for intermediate operations. This was the workflow used by Larry Wos to prove the challenge problem LCL073-1 with OTTER [83]. Our SGCD system also supports this, which played a major role in its ability to prove the aforementioned challenge problem.

Lemmas play a quite different role for a family of provers which we call *CM-CT* for *Connection Method/Clausal Tableaux*, exemplified by PTP [60], SETHEO [31], and leanCoP [45,44]. Underlying conceptual models are model elimination [33], clausal tableaux [29] and the CM. They enumerate proof structures while propagating variable bindings initialized by the goal through unification, and hence proceed in an inherently goal-driven way. While they are good at problems that benefit from goal direction, in general they are much weaker than RS provers. This is attributed to the fact that they do not re-use the proof of one subgoal as the solution of another: they do not use lemmas *internally*.

Lack of lemmas was identified early as a weakness of CM-CT, so there have been various proposed remedies [12,2,61,59,30,14,44,17]. Despite some insight and success, this did not yet elevate CM-CT to the level of the best RS systems. Nevertheless, the expectation remains that CM-CT provers would benefit from supplying lemmas as additional input. Hence, we included two CM-CT systems in our experiments, leanCoP and CMProver [9,70,71]. Two other systems consid-

ered here, SGCD and CCS [73], can be viewed as CM-CT systems extended to support specific forms of lemma generation and application.

Lemmas can be maintained within the prover as an inherent part of the method, as in saturation. They may also be created and applied by different systems, or different instances of the same system [10,53]. Larry Wos calls this *lemma adjunction* [82]. Lemmas created by one system are passed to a second system in two principal ways. First, they can be passed as *additional axioms*, in the hope that the second system finds a shorter proof in the wider but shallower search space. Second, external lemmas can be used to *replace search*. The second system then starts with the given lemmas as if they were the cached result of its previous computation. Moreover, the provided lemmas can be restricted in advance by heuristic methods, such as by a machine-learned model. SGCD supports this *replacing* lemma incorporation. The basic distinction between augmenting and replacing search with lemmas was already observed by Owen L. Astrachan and Mark E. Stickel [2] in the context of improving CM-CT provers.

2.1 Machine Learning for ATP

The past decade has seen numerous attempts to leverage machine learning in the automated theorem proving effort. Early systems mostly focused on premise selection, e.g. [67,1,69], aiming to reduce the number of axioms supplied as input to the prover. Other works provide internal guidance directly at the level of inferences during search, e.g. [32,23,15,25,84,52]. The emergence of generative language models has also led to some initial attempts at directly generating next proof steps, e.g. [47,66,48], moving the emphasis away from search.

In contrast to these lines of work, our focus is on learning the utility of lemmas. Close to our aims is [24,26], trying to identify globally useful lemmas in a collection of millions of proofs in HOL Light. Besides differences in the formal system, what distinguishes our work is that we learn a much more focused model: we put great emphasis on evaluating lemmas in the context of a particular goal and axiom set; in fact, our entire system was designed around the question whether a given lemma is moving the goal closer to the axioms. We argue that the D-term representation of all involved components (goal, lemma, axioms, proof) makes our framework particularly suitable for the lemma selection task.

We employ an iterative improvement approach first used in MaLAREa [67]: in each iteration, we run proof search guided by a learned model, extract training data from proofs, and fit a new model to the new data. These steps can be repeated profitably until performance saturates.

2.2 Condensed Detachment: Proofs as Terms

Condensed detachment (CD) was developed in the mid-1950s by Carew A. Meredith as an evolution of *substitution and detachment* [49,28,50,41]. Reasoning steps are by *detachment*, or modus ponens, under implicit substitution by most general unifiers. Its primary application is the investigation of axiomatizations of propositional logics at a first-order meta-level. CD also provides a technical

approach to the Curry-Howard correspondence, “formulas as types” [20,19] and is considered in witness theory [56]. Many early successes in ATP were on CD problems [38,65], but success was also found in the reverse direction. Refinements of the OTTER prover in the 1990s, some of which have found their ways into modern RS provers, were originally conceived and explored in the setting of CD [78,79,38,80,81,68,13,83].

From a first-order ATP perspective, a CD problem consists of *axioms*, i.e. positive unit clauses; a *goal theorem*, i.e. a single negative ground unit clause representing a universally-quantified atomic goal theorem after Skolemization; and the following ternary Horn clause that models detachment.

$$Det \stackrel{\text{def}}{=} P(i(x, y)) \wedge P(x) \rightarrow P(y).$$

The premises of *Det* are called the *major* and *minor* premise respectively. All atoms in the problem have the same predicate P , which is unary and stands for something like *provable*. The formulas of the investigated propositional logic are expressed as terms, where the binary function symbol i stands for *implies*.

CD may be seen as an *inference rule*. From an ATP perspective, a *CD inference step* can be described as a hyperresolution from *Det* and two positive unit clauses to a third positive unit clause. A *CD proof* is a proof of a CD problem constructed with the CD inference rule. CD proofs can be contrasted with other types of proof, such as a proof with binary resolution steps yielding non-unit clauses. Prover9 [36] chooses positive hyperresolution by default as its only inference rule for CD problems and thus produces CD proofs for these.

It is, however, another aspect of CD that makes it of particular interest for developing new ATP methods, which only recently came to our attention in the ATP context [74]: the structure of CD proofs can be represented in a very simple and convenient way as full binary trees, or as terms. In ATP we find this aspect in the CM, where the proof structure as a whole is in focus, in contrast to extending a set of formulas by deduction [8]. This view of CD is made precise and elaborated upon in [75], on which the subsequent informal presentation is based. We call the structure representations of CD proofs *D-terms*. A D-term is a term recursively built from numeral constants and the binary function symbol D whose arguments are D-terms. In other words, it is a full binary tree where the leaf nodes are labeled with constants. Four examples of D-terms are

$$1, \quad 2, \quad D(1, 1), \quad D(D(2, 1), D(1, D(2, 1))).$$

A D-term represents the structure of a proof. A proof in full is represented by a D-term together with a mapping of constant D-terms to axioms. Conversion between CD proofs and D-terms is straightforward: the use of an axiom corresponds to a constant D-term, while an inference step corresponds to a D-term $D(d_1, d_2)$ where d_1 is the D-term that proves the major premise and d_2 the minor.

Through first-order unification, constrained by axioms for the leaf nodes and the requirements of *Det* for inner nodes, it is possible to obtain a most general formula proven by a D-term [75]. We call it the *most general theorem* (MGT) of the D-term with respect to the axioms, unique up to renaming of variables. For

a given axiom map, not all D-terms necessarily have an MGT: if unification fails, we say the D-term has no MGT. It is also possible that different D-terms have the same MGT, or that the MGT of one is subsumed by the MGT of another. A D-term is a proof of the problem if its MGT subsumes the goal theorem.

As an example, let the constant D-term 1 be mapped to $P(i(x, i(x, x)))$, known as *Mingle* [65]. Then, the MGT of the D-term 1 is just this axiom. The MGT of the D-term $D(1, 1)$ is $P(i(x, i(x, x)), i(x, i(x, x)))$, that is, after renaming of variables, $P(y)\sigma$ where σ is the most general unifier of the set of pairs

$$\{\{P(i(x, y)), P(i(x', i(x', x')))\}, \{P(x), P(i(x'', i(x'', x'')))\}\}.$$

D-terms, as full binary trees, facilitate characterizing and investigating structural properties of proofs. While, for a variety of reasons, it is far from obvious how to measure the size of proofs obtained from ATP systems in general, for D-terms there are at least three straightforward size measures:

- The *tree size* of a D-term is the number of its inner nodes.
- The *height* of a D-term is the length of the longest root-leaf path.
- The *compacted size* of a D-term is the number of distinct compound subterms, or, in other words, the number of inner nodes of its minimal DAG.

In the literature compacted size is also called *length*, height *level* and tree size *CDcount* [68]. The D-term $D(D(1, D(1, 1)), D(D(1, 1), 1))$, for example, has tree size 5, compacted size 4 and height 3. *Factor equations* provide a compact way of writing D-terms: distinct subproofs with multiple incoming edges in the DAG receive numeric labels, by which they are referenced. The D-term $D(D(1, 1), D(D(1, D(1, 1)), D(1, D(1, 1))))$, for example, can be written as

$$2 = D(1, 1), \quad 3 = D(1, 2), \quad 4 = D(2, D(3, 3)).$$

2.3 Condensed Detachment for ATP and Lemmas

From an ATP point of view, the D-term, considered as a concept or also as a data structure, provides access to a proof as a whole. This exposes properties of proofs that are not merely local to the preconditions and conclusion of an inference step, but spread across the whole proof. The role of the calculus, which proceeds by inference steps, shifts to that of an inductive structure *specification*, rather than a recipe for incrementally building the structure by an efficient prover. Moreover, D-terms as objects provide insight into *all* proofs, all structures. For example, growth rates of the number of binary trees for tree size, height and compacted size are well-known with entries in *The On-Line Encyclopedia of Integer Sequences* [43] and provide upper bounds for the number of proofs [75]. A practical consequence for ATP is the justification of proof structure enumeration techniques where each structure appears at most once.

CD proofs suggest and allow a specific form of lemmas, which we call *unit* or *subtree* lemmas, reflecting two views on them. As formulas, they are positive unit clauses, which can be re-used in different CD inference steps. In the structural view, they are subterms, or subtrees, of the overall D-term. If they occur

multiply there, they are factored in the minimal DAG of the overall D-term. Both views are linked in that the formula of a lemma is the MGT of its D-term. The *compact size* measure specified above takes into account the compression achieved by unit/subtree lemmas. From the perspective of proof structure compression methods, unit/subtree lemmas have the property that the compression target is unique, because each tree is represented by a unique minimal DAG. CM-CT provers do not support such lemmas, which is the main reason for their notorious weakness on CD problems.

CD problems are restricted: only positive unit clauses, one negative ground clause, and one ternary Horn clause. Equality is not explicitly considered. Nevertheless, core characteristics of general first-order ATP problems are present: first-order variables, at least one binary function symbol and cyclic predicate dependency. The generalization to arbitrary Horn problems is not difficult [73].

2.4 SGCD — Structure Generating Theorem Proving

SGCD (*Structure Generating Theorem Proving for Condensed Detachment*) [72] is the central system used in our experiments as prover as well as lemma generator. It realizes an approach to first-order theorem proving combining techniques known from the CM and RS that was not fully recognized before. It generalizes (for CD problems) bottom-up preprocessing for and with CM-CT provers [59] and hypertableaux [4]. SGCD works by enumeration of proof structures together with unification of associated formulas, which is also the core method of the CM-CT provers. Structures for which unification fails are excluded immediately and do not appear in the enumeration. Typically, each structure appears at most once in the enumeration.

Let the proof structures be D-terms. Partition the set of all D-terms according to some *level* such that those in a lower level are strict subterms of those in a higher level. Tree size or height are examples of such a level. Let

$$\text{enum_dterm_mgt_pairs}(+Level, ?DTerm, ?Formula)$$

be a Prolog⁶ predicate enumerating D-terms and corresponding MGTs at a certain level, with respect to some quietly assumed axioms. We say that the predicate generates these pairs in an *axiom-driven* way. If the predicate is invoked with the formula argument instantiated by a ground formula, it enumerates D-terms that prove the formula at the specified level. The predicate is then used *goal-driven*, like a CM-CT prover. Invoking it for increasing level values realizes iterative deepening. There are further instantiation possibilities: if only the D-term is instantiated and the level is that of the D-term, its MGT is computed. If both D-term and formula are instantiated, the predicate acts as verifier.

The implementation includes several *generators*, concrete variants of the `enum_dterm_mgt_pairs` predicate for specific level characterizations. SGCD maintains a cache of $\langle level, D-term, formula \rangle$ triples obtained with a generator, which

⁶ Prolog serves here as a suitable specification language.

```

C := ∅;
for l := 0 to maxLevel do
  for m := l to l + preAddMaxLevel do
    enum_dterm_mgt_pairs(m, d, g);
    throw proof_found(d)
  N := {⟨l, d, f⟩ | enum_dterm_mgt_pairs(l, d, f)};
  if N = ∅ then throw exhausted;
  C := merge_news_into_cache(N, C)

```

Fig. 1. The nested loops of the SGCD theorem proving method.

is in turn used by the generator to obtain solutions for subproblems in levels below the calling level.

This cache is highly configurable. In particular, the number of entries can be limited, where only the best triples according to specified criteria are kept. Typical criteria are height or size of the formula, a heuristic shared with RS provers. Subsumed entries can be deleted, another feature in common with RS. Novel criteria are also supported, some of which relate the formula to the goal. Most criteria are based on the formula component of the triples, the MGT. Due to rigid variables [21], MGTs are not usually available in CM-CT provers [75] and cannot be used as a basis for heuristics.

From the perspective of lemma application, the cache use by the SGCD generators realizes *replacement* as discussed in Sect. 2. Search at lower levels than the calling level is replaced by cache retrieval, and replacement can be subject to heuristic restriction. SGCD further maintains a set of *abandoned* $\langle \text{level}, D\text{-term}, \text{formula} \rangle$ triples, those that are generated but do not qualify for entering the cache or were removed from the cache. These are kept as a source for heuristic evaluation of other triples and for lemma generation.

For theorem proving, SGCD proceeds as shown in Fig. 1. Input parameter g is the goal formula, while parameters $maxLevel$ and $preAddMaxLevel$ are configurable. `enum_dterm_mgt_pairs` represents a particular generator that is also configurable. It enumerates argument bindings nondeterministically: if it succeeds in the inner loop, an exception returns the D-term d . C is the cache. The procedure `merge_news_into_cache(N, C)` merges newly generated $\langle \text{level}, D\text{-term}, \text{formula} \rangle$ triples N into the cache C . If $maxLevel$ is configured as 0, the method proceeds in purely goal-driven mode with the inner loop performing iterative deepening on the level m . Similarity to CM-CT provers can be shown empirically by comparing the sets of solved TPTP problems [72]. Generally successful configurations of $preAddMaxLevel$ typically have values 0–3.

3 Iterative Improvement via Learned Lemma Selection

We take a fixed set of challenge problems along with a base theorem prover and aim to improve the prover via learning to identify useful lemmas. We first run the base prover on the problems and extract all D-term proofs found. In each D-term, we identify subterms as potential lemmas which are assigned utility scores.

Our aim is to build a *utility model* that takes a lemma as input and outputs a predicted utility, such that the model prediction closely matches the assigned utilities and generalises to unseen lemmas. To achieve this, we fix some class of parametric functions and optimise parameters via gradient descent. A large set of candidate lemmas is generated via structure enumeration, as described in Section 2, then each candidate is evaluated with the model and the top k lemmas are retained. These lemmas are then added to the original problem as axioms and we rerun proof search with the base prover. The extended set of proofs found can be used iteratively to repeat the above procedure.

3.1 Training Data Extraction

Given a proof tree, any connected subgraph can be considered a lemma. We focus on lemma subtrees where all leaves are axioms and leave more general subgraphs for future work. SGCD can also provide lemmas that are *not* part of the proof for use as training data. We assign a score to each lemma. We experimented with various scoring functions, but here we describe the two most successful:

1. **u_reproof**: We add lemmas one-by-one to the axioms and observe the change in (Prolog) inference steps required by the base prover. We normalise to $[-1, 1]$, where -1 means that the problem could not be solved within the original inference limit and 1 is assigned to the lemma that yields the greatest speedup.
2. **lf_is_in_proof**: This is 1 when the lemma is present in the proof and 0 otherwise, producing a binary classification problem. This score is only available for SGCD, since other provers do not provide non-proof lemmas.

3.2 Model Training

We experiment with gradient-descent optimisation for two classes of functions: linear models and graph neural networks (GNNs).

Linear Model The linear model is based on 51 novel manually-identified features, described in Appendix B. For each feature f_i there is an associated weight parameter w_i to produce the final predicted utility

$$U(\vec{f}; \vec{w}) = \sum_i f_i w_i$$

Graph Neural Network Model The second, more involved model is a GNN. Describing this model is beyond the scope of this paper: see e.g. [57] for a gentle introduction. What is crucial for our purposes is that no manual feature extraction is involved, the model processes D-terms of involved formulas directly and learns to extract useful features during optimisation. As input the model is given a graph, losslessly encoding D-terms of the lemma to be evaluated, the conjecture and the axioms. We use a graph neural network with 8 convolution layers

of 128 channels arranged into 4 residual blocks [16], followed by a hidden dense layer of 1024 neurons and a final dense layer that produces a single utility output. Batch normalization [22] is applied before each convolutional layer, and the non-linearity throughout is a rectified linear unit [42]. The precise configuration was found by manual optimisation with respect to a holdout set.

3.3 Candidate Lemma Generation

Candidate lemma generation is done with respect to a given problem’s axiom and goal. We use SGCD configured with the goal as g (Fig. 1) and *preAddMaxLevel* set to 0. SGCD then proceeds in essence axiom-driven, generating lemmas level by level. However, it does intersperse the goal-driven inner loop, which is only trying to prove the goal on the level directly above the last cached level. SGCD may terminate with a proof, in which case further lemma generation is pointless. Otherwise it terminates after *maxLevel* is reached, generation of new levels is exhausted, or a time limit is reached. We then use the cache C and the abandoned triples as the generated output lemmas.

Within this, there are many ways to configure SGCD. We obtained the best results generating by tree size and by PSP-level (explained below), combined with known-good heuristic restrictions. In particular we restrict the size of the lemma formulas to the maximum of the size of the axioms and the goal, multiplied by some factor (usually 2–5). We also restrict the number of elements in the cache, typically to 1,000. The lemmas are sorted by formula size measures, smaller preferred, to determine which are retained in the cache.

Proof structure generation by PSP-level is a novel technique based on an observation in proofs by Łukasiewicz and Meredith [72,75]. In a detachment step, often the D-term that proves one premise is a subterm of the term that proves the other. We turned this relationship into a proof structure enumeration method: structures in level $n+1$ are D-terms where one argument term is at level n and the other argument is a subterm of that term. The method is incomplete, but combines features of DAG enumeration while being compatible with a simple global lemma maintenance as realized with SGCD’s cache [75].

3.4 Lemma Selection and Iterative Improvement

We evaluate candidates with the learned model to select the top k lemmas for the next iteration of problem solving. It is then up to the prover to treat them as additional axioms, applicable to any prover, or to apply specialised treatment: in our experiments, provided lemmas are used to replace the inner lemma enumeration of SGCD and CCS-Vanilla. Before the obtained input lemmas are passed to a prover we supplement them with the lemmas for all their subproofs, i.e. we close the set of D-terms under the subterm relationship.

As long as there are new problems solved, we can iterate the entire process until performance saturates and no more proofs are added to the training dataset.

Table 1. Features of the considered provers: whether their proofs are available as D-terms (possibly after some conversion), whether they were used with *replacing* lemma incorporation (Sect. 2), whether they operate goal-driven, and the underlying method.

	SGCD	Prover9	CMProver	leanCoP	CCS-Vanilla	Vampire	E
D-terms	•	•	•	–	•	–	–
Repl. lemmas	•	–	–	–	•	–	–
Goal-driven	•/–	–	•	•	•	–	–
CM-CT	–	–	•	•	–	–	–
RS	–	•	–	–	–	•	•

We find that most provers plateau after 2–3 iterations, but some continue improving even after 5 iterations, such as some SGCD variants. To standardise presented experiments, we show up to two iterations.

3.5 Experiments with Iterative Improvement

Experiments described below show that several base provers improve with the addition of lemmas, and even that lemmas extracted with one prover are transferable to other theorem provers. Our experimental problem corpus is based on TPTP 8.1.2 [63]. We take 196 pure⁷ CD problems in the LCL domain, and enrich this set with single-axiom versions of all the problems to which a technique by Tarski [35], as specified by Rezuş [55], was applicable. This gives us 312 problems.

Table 1 gives an overview of the considered provers. If they produce D-terms as proofs, we have used these as training data. In the case of SGCD, after it has found a proof we also have access to internal lemmas, the cache and abandoned lemmas, from which we randomly select those that do not occur in the proof as negative training examples. Lemma incorporation by replacement in SGCD has many possibilities for configuration, but so far we used plausible settings determined automatically from properties of supplied lemmas. CCS-Vanilla is CCS [73] in a restricted configuration to find only those CD proofs with minimal compacted size, identifying problems that can clearly be solved with exhaustive search. It operates goal-driven, like the CM-CT provers, but by enumerating DAGs instead of trees through a local lemma maintenance mechanism. Vampire, for two decades the winner of CASC [62] in the FOF division, and its close rival E are RS provers that represent the state-of-the-art in first-order ATP.

All provers are recent public versions: Vampire 4.5.1, E 2.6, leanCoP 2.1. We provide results in terms of *time* limits, although for the Prolog provers SGCD, CMProver and CCS-Vanilla we used a roughly-equivalent inference limit to avoid fluctuations due to server workload. Our experiments showed virtually no difference between the two utility measures: `u_reproof` and `lf_is_in_proof`. Since `lf_is_in_proof` is faster to compute, but it is only applicable for SGCD, we use `lf_is_in_proof` for SGCD and `u_reproof` for all other provers.

⁷ TPTP contains 10 further CD problems which are satisfiable, have a non-atomic goal theorem, or use detachment with negation and disjunction.

Improving the Base Prover. Since our learning infrastructure assumes D-terms as proofs, our investigation is primarily directed towards base provers that emit proofs in this format, such as *SGCD*, *Prover9*, *CMProver* or *CCS-Vanilla*. Table 2 shows problems solved by these provers with two iterations of lemma selection, using the 200 best lemmas according to a linear utility model. The *Total* row gives the number of theorems proved by any of the three iterations shown. The stronger the base model, the harder it is to improve. *CMProver* and *CCS-Vanilla* are purely goal-driven and benefit greatly, reaching over 37% improvement for larger time limits. *SGCD* and *Prover9* improve over 5% for shorter time limits, but this effect gradually vanishes as the time limit is increased.

Table 2. Performance of different provers over 2 iterations of training a linear model.

Time	SGCD				Prover9				CMProver				CCS-Vanilla			
	50s	100s	500s	30m	50s	100s	500s	30m	50s	100s	500s	30m	50s	100s	500s	30m
Base	266	275	285	285	240	252	259	262	82	85	94	103	81	88	99	105
Iter 1	280	282	284	281	250	254	262	257	83	93	105	121	96	101	117	130
Iter 2	281	283	281	283	247	247	267	265	79	98	95	126	96	97	120	128
Total	282	284	286	286	253	258	269	267	91	105	112	141	106	105	133	145

Learned Lemmas to Enhance other Provers. While our learning setup is applicable only to provers that produce D-term proofs, any lemma identified as useful can be directly used by any other prover. In the following experiment we take a linear utility model trained with *SGCD* and use it to select 200 lemmas for *Vampire*, *E*, *Prover9* and *leanCoP*. Table 3 shows the greatest boost is for the purely goal-driven *leanCoP*, where there is over 40% improvement for all time limits. Second is *Vampire* with 8 – 15% improvement, followed by *Prover9* and *E* with around 3% improvement. Interestingly, *E* does not solve more problems with the lemmas, but it solves different ones, hence the improvement. These results suggest a great deal of transferability of the benefits of the lemma selector.

Table 3. Number of problems solved by *Vampire* (casc), *E* (autoschedule), *Prover9* and *leanCoP* without and with additional lemmas using various time limits.

Time	Vampire				E				Prover9				leanCoP			
	50s	100s	500s	30m	50s	100s	500s	30m	50s	100s	500s	30m	50s	100s	500s	30m
Base	221	224	252	263	253	264	275	281	236	244	257	260	70	71	77	77
Lemmas	249	257	274	283	256	266	275	275	246	250	261	269	100	103	111	113
Total	249	257	276	284	269	276	287	286	248	252	264	269	100	103	111	113

Changing the Number of Lemmas Added. Adding lemmas has potential to shorten proofs, but it also widens the search space, so it is not obvious how many lemmas are beneficial. Table 4 evaluates *SGCD* and *Vampire* with different numbers of lemmas. As little as 25 added lemmas yields substantial improvement, 7% for *Vampire* and 4% for *SGCD*, and performance does not drop as we add more lemmas: even at 500 we see no negative effect of the expanded search space.

Table 4. Number of problems solved by Vampire (casc) and SGCD as we alter the number of extracted lemmas. We use a time limit of 100s.

Lemma count	Vampire						SGCD					
	10	25	50	100	200	500	10	25	50	100	200	500
Base	227	227	227	227	227	227	275	275	275	275	275	275
Lemmas	226	242	246	258	257	258	278	285	284	281	283	284
Total	231	243	247	258	257	258	282	285	284	283	284	285

Linear vs GNN Model. The preceding subsections suggest that even a simple linear model can provide useful guidance when features are carefully selected. In Table 5 shows that the GNN — which processes the formulas directly and has no access to expert designed features — successfully learns to identify useful lemmas for SGCD and even slightly surpasses the linear model. LCL125-1 can only be solved by the GNN-assisted prover, even at extremely large time limits.

Table 5. Performance of SGCD over 2 iterations of training both a linear and a graph neural network model, for time limits 50 s, 100 s, 500 s and 30 min.

Time	Linear				GNN			
	50s	100s	500s	30m	50s	100s	500s	30m
Base	266	275	285	285	266	275	285	285
Iter 1	280	282	284	281	272	282	283	284
Iter 2	281	283	281	283	279	282	282	284
Total	282	284	286	286	279	285	287	287

4 Proving LCL073-1

LCL073-1 was proven by Meredith in the early 1950s with substitution and detachment [40] but it remains outstandingly hard for ATP, where it came to attention in 1992 [38]; TPTP reports rating 1.0 and status *Unknown* since 1997. Only Wos proved it in the year 2000 with several invocations of OTTER [83], transferring output and insight between runs. The problem has a single axiom,

$$P(i(i(i(i(x, y), i(n(z), n(u))), z), v), i(i(v, x), i(u, x)))),$$

and the goal $P(i(i(a, b), i(i(b, c), i(a, c))))$, known as *Syll* [65]. The wider context is showing that a single axiom entails the elements of a known axiomatization of a propositional logic. LCL038-1, discussed in [75], shows *Syll* from another single axiom. Experiments with SGCD in our workflow led to a proof of LCL073-1 (Fig. 2, also Appendix A) surprisingly quickly. Its compacted size is 46, between that of Meredith (40, reconstructed with CD in [83]) and that of Wos (74). Our workflow is much simpler than Wos', basically the same as our other experiments but restricted to one phase of lemma generation and incorporation, with only

PSP steps, i.e. one premise is a subproof of the other. It appears Wos’ “turning up the heat” [83] is now captured by PSP-level enumeration, a systematic proof enumeration technique resembling calculi for the CM rather than resolution.

5 Conclusion

We presented encouraging results about the use of lemmas in proof search. Provers are provided with lemmas generated via structure enumeration, a feature of the CM, and filtered with either learned guidance or manual heuristics. As a first step with this new methodology, we focus on the class of CD problems where we obtained remarkable success, strong results with our own system and substantial improvement of general first-order provers based on different paradigms, including the long-time competition leader *Vampire*. Moreover, our approach has led to the — in a sense first — automatic proof for the well-known Meredith single axiom problem with TPTP difficulty rating 1.0.

An important and novel aspect in our work was the explicit consideration of proof structures, which for CD have a particularly simple form in D-terms. Proof structures of the CM have a direct correspondence to these [75], such that the CM may guide the way to generalizations for more expressive logics. Another course of generalization is to move from unit lemmas, i.e. sharing of *subtrees* of D-terms, to more powerful lemmas. Preliminary work shows a correspondence between Horn clause lemmas, D-terms with variables, proofs in the connection structure calculus [12], and combinatory compression [73].

The learning-based experiments show little difference in performance between using a simple linear model and a more sophisticated graph neural network. We believe this is due to the small problem corpus, which yields a limited training signal. Hence, we plan to scale the system up to larger problem sets.

SGCD and its proof of LCL073-1 require, deserve and facilitate further systematic investigation. SGCD addresses a fundamental dilemma of modern ATP: calculi are too coarse and implementations too specific to adequately represent capabilities of provers. SGCD takes an algorithm-oriented view: systematic enumeration together with unification, caching and heuristic restrictions.

Our work also sheds new light on perspectives for the CM. It is well-known that the lack of inherent lemma maintenance is a disadvantage of the CM compared to resolution, which can be overcome with the connection structure calculus [12], a generalization of the CM. Here we see in experiments a drastic improvement of the CM-CT provers by supplementing their input with externally generated lemmas, however, they still lag behind RS provers. SGCD, which grew out of the CM-CT approach and integrates repeated lemma generation into the proving process, keeps up with RS provers on CD problems, and can even be applied to improve these by supplying its lemmas as additional input.

Acknowledgments. We thank Jens Otten for inspiring discussions at the outset of the current project.

References

1. Alemi, A.A., Chollet, F., Een, N., Irving, G., Szegedy, C., Urban, J.: Deepmath - Deep Sequence Models for Premise Selection. In: Lee, D., et al. (eds.) NIPS'16. pp. 2243–2251. Curran Associates Inc., USA (2016), <http://dl.acm.org/citation.cfm?id=3157096.3157347>
2. Astrachan, O.L., Stickel, M.E.: Caching and lemmaizing in model elimination theorem provers. In: Kapur, D. (ed.) CADE-11. pp. 224–238. Springer, Berlin (1992). https://doi.org/10.1007/3-540-55602-8_168
3. Bachmair, L., Ganzinger, H.: Resolution theorem proving. In: Robinson, A., Voronkov, A. (eds.) *Handb. of Autom. Reasoning*, vol. 1, chap. 2, pp. 19–99. Elsevier (2001). <https://doi.org/10.1016/B978-044450813-3/50004-7>
4. Baumgartner, P., Furbach, U., Niemelä, I.: Hyper tableaux. In: Alferes, J.J., Pereira, L.M., Orłowska, E. (eds.) JELIA'96. LNCS (LNAI), vol. 1126, pp. 1–17. Springer (1996). https://doi.org/10.1007/3-540-61630-6_1
5. Bibel, W.: *Automated Theorem Proving*. Vieweg, Braunschweig (1982). <https://doi.org/10.1007/978-3-322-90102-6>, second edition 1987
6. Bibel, W.: *Deduction: Automated Logic*. Academic Press, London (1993)
7. Bibel, W., Otten, J.: From Schütte's formal systems to modern automated deduction. In: Kahle, R., Rathjen, M. (eds.) *The Legacy of Kurt Schütte*, chap. 13, pp. 215–249. Springer (2020). https://doi.org/10.1007/978-3-030-49424-7_13
8. Bonacina, M.P.: A taxonomy of theorem-proving strategies. In: *Artificial Intelligence Today: Recent Trends and Developments*, pp. 43–84. Springer (2001)
9. Dahn, I., Wernhard, C.: First order proof problems extracted from an article in the Mizar mathematical library. In: Bonacina, M.P., Furbach, U. (eds.) FTP'97. pp. 58–62. RISC-Linz Report Series No. 97–50, Joh. Kepler Univ., Linz (1997), <https://www.logic.at/ftp97/papers/dahn.pdf>
10. Denzinger, J., Kronenburg, M., Schulz, S.: DISCOUNT — a distributed and learning equational prover. *J. Autom. Reasoning* **18**(2), 189 (1997)
11. Ebner, G., Hetzl, S., Leitsch, A., Reis, G., Weller, D.: On the generation of quantified lemmas. *J. Autom. Reasoning* **63**(1), 95–126 (2019). <https://doi.org/10.1007/s10817-018-9462-8>
12. Eder, E.: A comparison of the resolution calculus and the connection method, and a new calculus generalizing both methods. In: Börger, E., Kleine Büning, H., Richter, M.M. (eds.) CSL '88. LNCS, vol. 385, pp. 80–98. Springer (1989). <https://doi.org/10.1007/BFb0026296>
13. Fitelson, B., Wos, L.: Missing proofs found. *J. Autom. Reasoning* **27**(2), 201–225 (2001). <https://doi.org/10.1023/A:1010695827789>
14. Fuchs, M.: Lemma generation for model elimination by combining top-down and bottom-up inference. In: Dean, T. (ed.) IJCAI 1999. pp. 4–9. Morgan Kaufmann (1999), <http://ijcai.org/Proceedings/99-1/Papers/001.pdf>
15. Gauthier, T., Kaliszky, C., Urban, J., Kumar, R., Norrish, M.: Learning to prove with tactics. CoRR [abs/1804.00596](https://arxiv.org/abs/1804.00596) (2018), <http://arxiv.org/abs/1804.00596>
16. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR 2016. pp. 770–778 (2016)
17. Hester, J.: *Novel Methods for First Order Automated Theorem Proving*. Ph.D. thesis, University of Florida (2021)
18. Hetzl, S., Leitsch, A., Weller, D.: Towards algorithmic cut-introduction. In: Bjørner, N., Voronkov, A. (eds.) LPAR 2012. LNCS, vol. 7180, pp. 228–242. Springer (2012). https://doi.org/10.1007/978-3-642-28717-6_19

19. Hindley, J.R.: Basic Simple Type Theory. Cambridge University Press (1997). <https://doi.org/10.1017/CBO9780511608865>
20. Hindley, J.R., Meredith, D.: Principal type-schemes and condensed detachment. *Journal of Symbolic Logic* **55**(1), 90–105 (1990). <https://doi.org/10.2307/2274956>
21. Hähnle, R.: Tableaux and related methods. In: Robinson, A., Voronkov, A. (eds.) *Handb. of Autom. Reasoning*, vol. 1, chap. 3, pp. 101–178. Elsevier (2001). <https://doi.org/10.1016/b978-044450813-3/50005-9>
22. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: *ICML 2015*. pp. 448–456. *Proceedings of Machine Learning Research* (2015)
23. Jakubuv, J., Urban, J.: ENIGMA: efficient learning-based inference guiding machine. In: Geuvers, H., England, M., Hasan, O., Rabe, F., Teschke, O. (eds.) *CICM 2017*. LNCS, vol. 10383, pp. 292–302. Springer (2017). https://doi.org/10.1007/978-3-319-62075-6_20, https://doi.org/10.1007/978-3-319-62075-6_20
24. Kaliszyk, C., Urban, J.: Learning-assisted theorem proving with millions of lemmas. *Journal of Symbolic Computation* **69**, 109–128 (2015). <https://doi.org/https://doi.org/10.1016/j.jsc.2014.09.032>, <https://www.sciencedirect.com/science/article/pii/S074771711400100X>, *symbolic Computation in Software Science*
25. Kaliszyk, C., Urban, J., Michalewski, H., Olsák, M.: Reinforcement learning of theorem proving. In: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (eds.) *NeurIPS 2018*. pp. 8836–8847 (2018), <https://papers.nips.cc/paper/2018/file/55acf8539596d25624059980986aaa78-Paper.pdf>
26. Kaliszyk, C., Urban, J., Vyskocil, J.: Lemmatization for stronger reasoning in large theories. In: Lutz, C., Ranise, S. (eds.) *Frontiers of Combining Systems - 10th International Symposium, FroCoS 2015, Wroclaw, Poland, September 21-24, 2015*. *Proceedings*. LNCS, vol. 9322, pp. 341–356. Springer (2015). https://doi.org/10.1007/978-3-319-24246-0_21, https://doi.org/10.1007/978-3-319-24246-0_21
27. Kovács, L., Voronkov, A.: First-order theorem proving and VAMPIRE. In: *Computer Aided Verification: 25th International Conference, CAV 2013, Saint Petersburg, Russia, July 13-19, 2013*. *Proceedings* 25. pp. 1–35. Springer (2013)
28. Lemmon, E.J., Meredith, C.A., Meredith, D., Prior, A.N., Thomas, I.: *Calculi of pure strict implication*. In: Davis, J.W., Hockney, D.J., Wilson, W.K. (eds.) *Philosophical Logic*, pp. 215–250. Springer Netherlands, Dordrecht (1969). https://doi.org/10.1007/978-94-010-9614-0_17, reprint of a technical report, Canterbury University College, Christchurch, 1957
29. Letz, R.: *Tableau and Connection Calculi. Structure, Complexity, Implementation*. *Habilitationsschrift*, TU München (1999), available from <http://www2.tcs.ifi.lmu.de/~letz/habil.ps>, accessed Jun 30, 2022
30. Letz, R., Mayr, K., Goller, C.: Controlled integration of the cut rule into connection tableaux calculi. *J. Autom. Reasoning* **13**(3), 297–337 (1994)
31. Letz, R., Schumann, J., Bayerl, S., Bibel, W.: SETHEO: A high-performance theorem prover. *J. Autom. Reasoning* **8**(2), 183–212 (1992). <https://doi.org/10.1007/BF00244282>
32. Loos, S.M., Irving, G., Szegedy, C., Kaliszyk, C.: Deep network guided proof search. In: Eiter, T., Sands, D. (eds.) *LPAR-21. EPIc*, vol. 56, pp. 85–105 (2017). <https://doi.org/10.29007/8mwc>
33. Loveland, D.W.: *Automated Theorem Proving: A Logical Basis*. North-Holland, Amsterdam (1978)
34. Lukasiewicz, J.: *Selected Works*. North Holland (1970), edited by L. Borkowski

35. Łukasiewicz, J., Tarski, A.: Untersuchungen über den Aussagenkalkül. *Comptes rendus des séances de la Soc. d. Sciences et d. Lettres de Varsovie* **23** (1930), English translation in [34], p. 131–152
36. McCune, W.: Prover9 and Mace4 (2005–2010), <http://www.cs.unm.edu/~mccune/prover9>
37. McCune, W.: OTTER 3.3 Reference Manual. Tech. Rep. ANL/MCS-TM-263, Argonne National Laboratory (2003), <https://www.cs.unm.edu/~mccune/otter/Otter33.pdf>, accessed Jun 30, 2022
38. McCune, W., Wos, L.: Experiments in automated deduction with condensed detachment. In: Kapur, D. (ed.) CADE-11. LNCS (LNAI), vol. 607, pp. 209–223. Springer (1992). https://doi.org/10.1007/3-540-55602-8_167
39. Meredith, C.A., Prior, A.N.: Notes on the axiomatics of the propositional calculus. *Notre Dame J. of Formal Logic* **4**(3), 171–187 (1963). <https://doi.org/10.1305/ndjfl/1093957574>
40. Meredith, C.A.: Single axioms for the systems (C, N), (C, O) and (A, N) of the two-valued propositional calculus. *J. Computing Systems* **1**, 155–164 (1953)
41. Meredith, D.: In memoriam: Carew Arthur Meredith (1904–1976). *Notre Dame J. of Formal Logic* **18**(4), 513–516 (Oct 1977). <https://doi.org/10.1305/ndjfl/1093888116>
42. Nair, V., Hinton, G.E.: Rectified linear units improve restricted Boltzmann machines. In: ICML 2010. pp. 807–814 (2010)
43. OEIS Foundation Inc.: The On-Line Encyclopedia of Integer Sequences (2021), <http://oeis.org>
44. Otten, J.: Restricting backtracking in connection calculi. *AI Communications* **23**(2-3), 159–182 (2010). <https://doi.org/10.3233/AIC-2010-0464>
45. Otten, J., Bibel, W.: leanCoP: lean connection-based theorem proving. *J. Symb. Comput.* **36**(1-2), 139–161 (2003). [https://doi.org/10.1016/S0747-7171\(03\)00037-3](https://doi.org/10.1016/S0747-7171(03)00037-3)
46. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: *Advances in Neural Information Processing Systems* **32**, pp. 8024–8035. Curran Associates, Inc. (2019), <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
47. Piotrowski, B., Urban, J.: Guiding inferences in connection tableau by recurrent neural networks. In: Benz Müller, C., Miller, B.R. (eds.) CICM. LNCS, vol. 12236, pp. 309–314. Springer (2020). https://doi.org/10.1007/978-3-030-53518-6_23, https://doi.org/10.1007/978-3-030-53518-6_23
48. Polu, S., Sutskever, I.: Generative language modeling for automated theorem proving. *CoRR* **abs/2009.03393** (2020), <https://arxiv.org/abs/2009.03393>
49. Prior, A.N.: Logicians at play; or Syll, Simp and Hilbert. *Australasian Journal of Philosophy* **34**(3), 182–192 (1956). <https://doi.org/10.1080/00048405685200181>
50. Prior, A.N.: *Formal Logic*. Clarendon Press, Oxford, 2nd edn. (1962). <https://doi.org/10.1093/acprof:oso/9780198241560.001.0001>
51. Pudlák, P.: Search for faster and shorter proofs using machine generated lemmas. In: Sutcliffe, G., Schmidt, R., Schulz, S. (eds.) ESCoR 2006. CEUR Workshop Proc., vol. 192, pp. 34–53. CEUR-WS.org (2006), <http://ceur-ws.org/Vol-192/paper03.pdf>

52. Rawson, M., Reger, G.: lazyCoP: Lazy paramodulation meets neurally guided search. In: Das, A., Negri, S. (eds.) TABLEAUX 2021. LNCS (LNAI), vol. 12842, pp. 187–199. Springer (2021). https://doi.org/10.1007/978-3-030-86059-2_11
53. Reger, G., Tishkovsky, D., Voronkov, A.: Cooperating proof attempts. In: CADE-25. pp. 339–355. Springer (2015). https://doi.org/10.1007/978-3-319-21401-6_23
54. Rezuş, A.: The Curry-Howard Correspondence revisited (2019b). In: Witness Theory – Notes on λ -calculus and Logic [56], pp. 209–214
55. Rezuş, A.: Tarski’s Claim thirty years later (2010). In: Witness Theory – Notes on λ -calculus and Logic [56], pp. 217–225, preprint (2016): <http://www.equivalences.org/editions/proof-theory/ar-tc-20160512.pdf>
56. Rezuş, A.: Witness Theory – Notes on λ -calculus and Logic, Studies in Logic, vol. 84. College Publications, London (2020)
57. Sanchez-Lengeling, B., Reif, E., Pearce, A., Wiltshcko, A.B.: A gentle introduction to graph neural networks. Distill (2021). <https://doi.org/10.23915/distill.00033>, <https://distill.pub/2021/gnn-intro>
58. Schulz, S., Cruanes, S., Vukmirović, P.: Faster, higher, stronger: E 2.3. In: Fontaine, P. (ed.) CADE 27. pp. 495–507. No. 11716 in LNAI, Springer (2019). https://doi.org/10.1007/978-3-030-29436-6_29
59. Schumann, J.M.P.: DELTA – A bottom-up preprocessor for top-down theorem provers. In: CADE-12. LNCS (LNAI), vol. 814, pp. 774–777. Springer (1994). https://doi.org/10.1007/3-540-58156-1_58
60. Stickel, M.E.: A Prolog technology theorem prover: implementation by an extended Prolog compiler. J. Autom. Reasoning **4**(4), 353–380 (1988). <https://doi.org/10.1007/BF00297245>
61. Stickel, M.E.: Upside-down meta-interpretation of the model elimination theorem-proving procedure for deduction and abduction. J. Autom. Reasoning **13**(2), 189–210 (1994). <https://doi.org/10.1007/BF00881955>
62. Sutcliffe, G.: The CADE ATP system competition — CASC. AI Magazine **37**(2), 99–101 (2016)
63. Sutcliffe, G.: The TPTP problem library and associated infrastructure. From CNF to TH0, TPTP v6.4.0. J. Autom. Reasoning **59**(4), 483–502 (2017)
64. Sutcliffe, G., Gao, Y., Colton, S.: A grand challenge of theorem discovery. In: Worksh. Challenges and Novel Applications for Automated Reasoning, 19th IJ-CAR. pp. 1–11 (2003), online: https://www.cs.miami.edu/home/geoff/Papers/Conference/2003_SGC03_CNAAR-1-11.pdf
65. Ulrich, D.: A legacy recalled and a tradition continued. J. Autom. Reasoning **27**(2), 97–122 (2001). <https://doi.org/10.1023/A:1010683508225>
66. Urban, J., Jakubuv, J.: First neural conjecturing datasets and experiments. In: Benz Müller, C., Miller, B.R. (eds.) CICM 2020. LNCS, vol. 12236, pp. 315–323. Springer (2020). https://doi.org/10.1007/978-3-030-53518-6_24, https://doi.org/10.1007/978-3-030-53518-6_24
67. Urban, J., Sutcliffe, G., Pudlák, P., Vyskocil, J.: MaLAREa SG1 – Machine Learner for Automated Reasoning with Semantic Guidance. In: Armando, A., Baumgartner, P., Dowek, G. (eds.) IJCAR 2008. LNCS, vol. 5195, pp. 441–456. Springer (2008). https://doi.org/10.1007/978-3-540-71070-7_37
68. Veroff, R.: Finding shortest proofs: An application of linked inference rules. J. Autom. Reasoning **27**(2), 123–139 (2001). <https://doi.org/10.1023/A:1010635625063>
69. Wang, M., Tang, Y., Wang, J., Deng, J.: Premise selection for theorem proving by deep graph embedding. In: Guyon, I., et al. (eds.) NIPS 2017. pp. 2783–2793 (2017), <http://papers.nips.cc/paper/6871-premise-selection-for-theorem-proving-by-deep-graph-embedding>

70. Wernhard, C.: The PIE system for proving, interpolating and eliminating. In: Fontaine, P., Schulz, S., Urban, J. (eds.) PAAR 2016. CEUR Workshop Proc., vol. 1635, pp. 125–138. CEUR-WS.org (2016), <http://ceur-ws.org/Vol-1635/paper-11.pdf>
71. Wernhard, C.: Facets of the PIE environment for proving, interpolating and eliminating on the basis of first-order logic. In: Hofstedt, P., et al. (eds.) DECLARE 2019. LNCS (LNAI), vol. 12057, pp. 160–177 (2020). https://doi.org/10.1007/978-3-030-46714-2_11
72. Wernhard, C.: CD Tools – Condensed detachment and structure generating theorem proving (system description) (2022), <https://arxiv.org/abs/2207.08453>
73. Wernhard, C.: Generating compressed combinatory proof structures – an approach to automated first-order theorem proving. In: Konev, B., Schon, C., Steen, A. (eds.) PAAR 2022. CEUR Workshop Proc., vol. 3201. CEUR-WS.org (2022), <https://arxiv.org/abs/2209.12592>
74. Wernhard, C., Bibel, W.: Learning from Łukasiewicz and Meredith: Investigations into proof structures. In: Platzer, A., Sutcliffe, G. (eds.) CADE 28. LNCS (LNAI), vol. 12699, pp. 58–75. Springer (2021). https://doi.org/10.1007/978-3-030-79876-5_4
75. Wernhard, C., Bibel, W.: Investigations into proof structures (2023), preprint, <http://cs.christophwernhard.com/papers/investigations/>
76. Wielemaker, J., Schrijvers, T., Triska, M., Lager, T.: SWI-Prolog. Theory and Practice of Logic Programming **12**(1-2), 67–96 (2012). <https://doi.org/10.1017/S1471068411000494>
77. Woltzenlogel Paleo, B.: Atomic cut introduction by resolution: Proof structuring and compression. In: Clarke, E.M., Voronkov, A. (eds.) LPAR-16. LNCS, vol. 6355, pp. 463–480. Springer (2010). https://doi.org/10.1007/978-3-642-17511-4_26
78. Wos, L., Winker, S., McCune, W., Overbeek, R., Lusk, E., Stevens, R., Butler, R.: Automated reasoning contributes to mathematics and logic. In: Stickel, M.E. (ed.) CADE-10. pp. 485–499. Springer (1990). https://doi.org/10.1007/3-540-52885-7_109
79. Wos, L.: Automated reasoning and Bledsoe’s dream for the field. In: Boyer, R.S. (ed.) Automated Reasoning: Essays in Honor of Woody Bledsoe, pp. 297–345. Automated Reasoning Series, Kluwer Academic Publishers (1991). https://doi.org/10.1007/978-94-011-3488-0_15
80. Wos, L.: The resonance strategy. Computers Math. Applic. **29**(2), 133–178 (1995). [https://doi.org/10.1016/0898-1221\(94\)00220-F](https://doi.org/10.1016/0898-1221(94)00220-F)
81. Wos, L.: The power of combining resonance with heat. J. Autom. Reasoning **17**(1), 23–81 (1996). <https://doi.org/10.1007/BF00247668>
82. Wos, L.: Lemma inclusion versus lemma adjunction. Association for Automated Reasoning Newsletter **44** (September 1999), online: <https://aarinc.org/Newsletters/044-1999-09.html>, accessed 24 Feb 2023
83. Wos, L.: Conquering the Meredith single axiom. J. Autom. Reasoning **27**(2), 175–199 (2001). <https://doi.org/10.1023/A:1010691726881>
84. Zombori, Z., Urban, J., Brown, C.E.: Prolog technology reinforcement learning prover - (system description). In: Peltier, N., Sofronie-Stokkermans, V. (eds.) IJ-CAR 2020. LNCS, vol. 12167, pp. 489–507. Springer (2020). https://doi.org/10.1007/978-3-030-51054-1_33, https://doi.org/10.1007/978-3-030-51054-1_33

A Our Proof of LCL073-1 with MGTs

1. $CCCCCpqCNrNsrtCCtpCsp$
2. $CCppqCrq = D1D1D11$
3. $CpCqCrr = D22$
4. $CCpCqqrCsr = D13$
5. $CCppqrCqr = D14$
6. $CpCCpqCrq = D51$
7. $CpCCCqqrCsr = D56$
8. $CCppqrCCCqsCNrNpr = DDD1D176n$
9. $CCppqCNCCCrpsCtsNrCCCrpsCts = D86$
10. $CCppqCNCrpNCCCsCptuCvuCrp = D8D19$
11. $CpCCqrCCNqNpr = DD1D1D4.10.1$
12. $CCpNqCCqrCsrCtCCqrCsr = D1D6D1DD1D9D9DD11.341$
13. $CCpCqrCCCsCtNprCqr = DDD12.D5D8.12.n7$
14. $CCCCCpqrstCCCqrst = D1D13.D1D13.5$
- *15. $CCpqCCqrCpr = DD1D13.DDDD13.69.11.10.D14.D14.1$

Fig. 3. Our proof of LCL073-1 in the notation of Meredith [39,75]. For each factor of the overall D-term the argument term of its MGT is there shown in Łukasiewicz's notation.

1. $P(i(i(i(i(p, q), i(n(r), n(s))), r), t), i(i(t, p), i(s, p))))$
2. $P(i(i(i(p, p), q), i(r, q))) = D(1, D(1, D(1, 1)))$
3. $P(i(p, i(q, i(r, r)))) = D(2, 2)$
4. $P(i(i(i(p, i(q, q)), r), i(s, r))) = D(1, 3)$
5. $P(i(i(i(p, q), r), i(q, r))) = D(1, 4)$
6. $P(i(p, i(i(p, q), i(r, q)))) = D(5, 1)$
7. $P(i(p, i(i(i(q, p), r), i(s, r)))) = D(5, 6)$
8. $P(i(i(i(i(p, q), r), i(i(i(q, s), i(n(r), n(p))), r)))) = D(D(D(1, D(1, 7)), 6), n)$
9. $P(i(i(i(p, q), i(n(i(i(r, p), s), i(t, s))), n(r))), i(i(i(r, p), s), i(t, s))) = D(8, 6)$
10. $P(i(i(i(p, q), i(n(i(r, p))), n(i(i(i(s, i(p, t)), u), i(v, u))))), i(r, p)) = D(8, D(1, 9))$
11. $P(i(p, i(i(q, r), i(i(n(q), n(p))), r))) = D(D(1, D(1, D(4, 10))), 1)$
12. $P(i(i(i(p, n(q)), i(i(q, r), i(s, r))), i(t, i(i(q, r), i(s, r))))))$
 $= D(1, D(6, D(1, D(D(1, D(9, D(9, D(D(11, 3), 4))))), 1)))$
13. $P(i(i(p, i(q, r)), i(i(i(s, i(t, n(p))), r), i(q, r)))) = D(D(D(12, D(5, D(8, 12))), n), 7)$
14. $P(i(i(i(i(p, q), r), s), t), i(i(i(q, r), s), t))) = D(1, D(13, D(1, D(13, 5))))$
- *15. $P(i(i(p, q), i(i(q, r), i(p, r))))$
 $= D(D(1, D(13, D(D(D(13, 6), 9), 11), 10))), D(14, D(14, 1)))$

Fig. 4. Our proof of LCL073-1 with the ATP-oriented notation for formulas and D-terms of the paper, arranged such that it mimicks Meredith's notation.

B Manually Extracted Lemma Features

Below we specify the features that are manually extracted from lemmas. We also provide their types as follows:

- *NatNum* A natural number.
- *NormalizedValue* A number between 0 and 1.

B.1 Features of the Lemma’s D-term

- lf_d_csize** : *NatNum*
 Compacted size, tree size and height of the lemma’s D-term.
- lf_d_tsize** : *NatNum*
 Compacted size of the lemma’s D-term after replacing all variables with 0.
- lf_d_height** : *NatNum*
 Compacted size, tree size and height of the lemma’s D-term.
- lf_d_grd_csize** : *NatNum*
 Compacted size of the lemma’s D-term after replacing all variables with 0.
- lf_d_major_minor_relation** : *NatNum*
 Describes the structural relationship of the subproofs of the major and minor premise of the lemma. Can have the following values: 0 identical or D-term is atomic; 1 is a strict superterm; 2 is a strict subterm; 3 none of these relationships; 4 for nonground D-terms.
- lf_d_number_of_terminals** : *NatNum*
 Number of subterms in the lemma’s D-term which are of the form $d(d_1, d_2)$ where neither of d_1, d_2 is a compound term.

B.2 Context Dependent Features of the Lemma’s D-Term

The features depend on the lemma as embedded in a given proof. They are only available for training data extracted from proofs.

- lfp_containing_proof** : **proof**
 The proof that contains the lemma. Specified as the atom that identifies the proof. See also **lf_proof**.
- lfp_d_occs** : *NatNum*
 If the lemma’s D-term is ground, the number of occurrences of the lemma’s D-term in the proof’s D-term. If the lemma’s D-term is non-ground, the number of subterm occurrences of the proof’s D-term that are instances of the lemma’s D-term (note that these subterm occurrences may overlap).
- lfp_d_incoming** : *NatNum*
 The number of incoming edges of the lemma in the minimal DAG representing the proof’s D-term. This can not be meaningfully determined for all lemma computation methods that yield s-lemmas.
- lfp_d_occs_innermost_matches** : *NatNum*
 If the lemma’s D-term is ground, the same as **lf_d_occs**. If the lemma’s D-term has variables, the number of subterm occurrences of the proof’s D-term that would be rewritten by INNERMOST replacement.

- `lfp_d_occs_outermost_matches` : *NatNum*
 If the lemma's D-term is ground, the same as `lf_d_occs`. If the lemma's D-term has variables, the number of subterm occurrences of the proof's D-term that would be rewritten by OUTERMOST replacement.
- `lfp_d_min_goal_dist` : *NatNum*
 Number of edges in the proof's D-term of the shortest downward path from the root to a subtree that is an instance of the lemma's D-term.

B.3 Special Features of the Lemma's Formula Components

- `lf_b_length` : *NatNum*
 Length of the *Body* component of the lemma's formula.
- `lf_hb_distinct_hb_shared_vars` : *NatNum*
 Number of distinct variables that occur in the *Head* as well as in the *Body* component of the lemma's formula.
- `lf_hb_distinct_h_only_vars` : *NatNum*
 Number of distinct variables that occur in the *Head* but not the *Body* component of the lemma's formula.
- `lf_hb_distinct_b_only_vars` : *NatNum*
 Number of distinct variables that occur in the *Body* but not the *Head* component of the lemma's formula.
- `lf_hb_singletons` : *NatNum*
 Number of distinct variables with a single occurrence in *Head* and *Body* taken together.
- `lf_hb_double_negation_occs` : *NatNum*
 Number of instances of $\mathbf{n}(\mathbf{n}(_))$ in *Head* and *Body*.
- `lf_hb_nongoal_symbol_occs` : *NatNum*
 Number of occurrences of symbols (functions, constants) in *Head* and *Body* that do not appear in the problem's goal. Context-dependent in that it refers to the problem's goal formula.
- `lf_h_excluded_goal_subterms` : *NatNum*
 Number of distinct subterms of the goal formula (after replacing constants systematicall by variables) that are not a subterm of *Head* (modulo the variant relationship). Context-dependent in that it refers to the problem's goal formula.
- `lf_h_subterms_not_in_goal` : *NatNum*
 Number of distinct subterms of the head that are not a subterm of the goal formula (after replacing constants systematicall by variables, modulo the variant relationship). Context-dependent in that it refers to the problem's goal formula.
- `lf_hb_compression_ratio_raw_deflate` : *NormalizedValue*
- `lf_hb_compression_ratio_treerepair` : *NormalizedValue*
- `lf_hb_compression_ratio_dag` : *NormalizedValue*

Indicates how much the lemma’s formula can be compressed. The value is roughly compressed size divided by original tree size. That is, formulas with “much regularity” such that they can be compressed stronger receive smaller values. The different properties realize this in variants for different notions and implementations of compression. The `raw_deflate` version depends on intrinsics of SWI-Prolog’s term representation and possibly gives different results for the same formula, depending on how it internally shares subterms.

`lf_hb_organic` : *NatNum*

Whether the formula is organic. A nonempty *Body* is translated to an implication, e.g., if *Body* = [a,b,c], the considered formula is `i(a,i(b,i(c,Head)))`. Determined with `Minisat`. Values: 0: the formula is organic; 1 the formula is not organic but weakly organic; 2 the formula is not weakly organic. See also http://cs.christophwernhard.com/cdtools/downloads/cdtools/pldoc/organic_cd.html.

`lf_hb_name` : *Atom*

A name of the formula if it is well known under some name. For a formula with nonempty body the translation to implication is considered (as for `lf_hb_organic`) and the name is prefixed with `meta_`. If the formula is not known under some name, the value is `zzz`. See also `named_axiom/2` in http://cs.christophwernhard.com/cdtools/downloads/cdtools/pldoc/named_axioms_cd.html.

`lf_hb_name_status` : *NatNum*

Number indicating whether the formula has a name in the sense of `lf_hb_name`: 0 if it has a name and an empty body, 1 if it has a nonempty body and a name (prefixed with `meta_`), 2 otherwise.

B.4 General Features of the Lemma’s Formula Components

These features are specified below schematically with *COMP* for `h`, `b`, and `hb`, referring the respective features for the *Head* component, the *Body* component and both of the components joined together. The schema parameter *ITEM* for `var`, `const` and `fun` refers to variables, constants (atomic values in Prolog syntax) and function symbols with arity ≥ 1 , respectively.

`lf_COMP_csize` : *NatNum*

`lf_COMP_tsize` : *NatNum*

`lf_COMP_height` : *NatNum*

Compacted size, tree size and height, respectively. (We use these notions, which are also used for D-terms, here for formula terms.)

`lf_COMP_distinct_vars` : *NatNum*

Number of distinct variables.

`lf_COMP_ITEM_occs` : *NatNum*

Number of occurrences of syntactic objects of kind *ITEM*.

`lf_COMP_occs_of_most_frequent_ITEM` : *NatNum*

Maximum number of occurrences of a syntactic object of kind *ITEM*.