



# Certified Knowledge Compilation with Application to Verified Model Counting

Randal E. Bryant  

Computer Science Department, Carnegie Mellon University, Pittsburgh, PA 15213 USA

Wojciech Nawrocki  

Department of Philosophy, Carnegie Mellon University

Jeremy Avigad  

Department of Philosophy, Carnegie Mellon University

Marijn J. H. Heule  

Computer Science Department, Carnegie Mellon University

---

## Abstract

Computing many useful properties of Boolean formulas, such as their weighted or unweighted model count, is intractable on general representations. It can become tractable when formulas are expressed in a special form, such as the deterministic, decomposable, negation normal form (d-DNNF). *Knowledge compilation* is the process of converting a formula into such a form. Unfortunately existing knowledge compilers provide no guarantee that their output correctly represents the original formula, and therefore they cannot validate a model count, or any other computed value.

We present *Partitioned-Operation Graphs* (POGs), a generalization of the representations used by existing knowledge compilers. We have designed CPOG, a framework that can express proofs of equivalence between a POG and a Boolean formula in conjunctive normal form (CNF).

We have developed a program that generates POG representations from the d-DNNF graphs produced by the state-of-the-art knowledge compiler D4, as well as checkable CPOG proofs certifying that the output POGs are equivalent to the input CNF formulas. Our toolchain for generating and verifying POGs scales to all but the largest d-DNNF graphs produced by D4 for formulas from a recent model counting competition. Additionally, we have developed a formally verified CPOG checker and model counter for POGs in the Lean 4 proof assistant. In doing so, we proved the soundness of our proof framework. These programs comprise the first formally verified toolchain for weighted and unweighted model counting.

**2012 ACM Subject Classification** Theory of computation → Automated reasoning

**Keywords and phrases** Propositional model counting, Proof checking

**Digital Object Identifier** 10.4230/LIPIcs...

**Funding** *Randal E. Bryant*: Supported by NSF grant CCF-2108521

*Wojciech Nawrocki*: Hoskinson Center for Formal Mathematics

*Jeremy Avigad*: Hoskinson Center for Formal Mathematics

*Marijn J. H. Heule*: Supported by NSF grant CCF-2108521

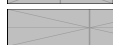
## 1 Introduction

Given a Boolean formula  $\phi$ , modern Boolean satisfiability (SAT) solvers can find an assignment satisfying  $\phi$  or generate a proof that no such assignment exists. They have applications across a variety of domains including computational mathematics, combinatorial optimization, and the formal verification of hardware, software, and security protocols. Some applications, however, require going beyond Boolean satisfiability. For example, the *model counting problem* requires computing the number of satisfying assignments of a formula, including in cases where there are far too many to enumerate individually. Model counting has applications in artificial intelligence, computer security, and statistical sampling. There are also many useful



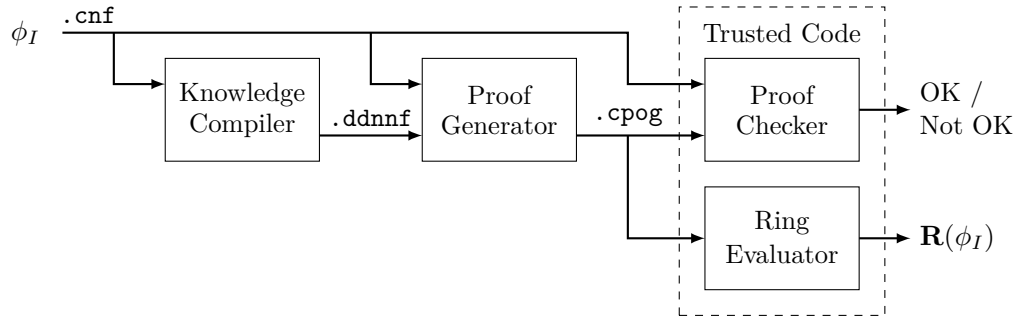
© Author: Please provide a copyright holder;

licensed under Creative Commons License CC-BY 4.0



Leibniz International Proceedings in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany



■ **Figure 1** Certifying toolchain. The output of a standard knowledge compiler is converted into a combined graph/proof (CPOG) which can be independently checked and evaluated

45 extensions of standard model counting, including *weighted model counting*, where a weight  
 46 is defined for each possible assignment, and the goal becomes to compute the sum of the  
 47 weights of the satisfying assignments.

48 Model counting is a challenging problem—more challenging than the already NP-hard  
 49 Boolean satisfiability. Several tractable variants of Boolean satisfiability, including 2-SAT,  
 50 become intractable when the goal is to count models and not just determine satisfiability  
 51 [21]. Nonetheless, a number of model counters that scale to very large formulas have been  
 52 developed, as witnessed by the progress in recent model counting competitions.

53 One approach to model counting, known as *knowledge compilation*, transforms the  
 54 formula into a structured form for which model counting is straightforward. For example, the  
 55 *deterministic, decomposable, negation normal form* (d-DNNF) introduced by Darwiche [4, 5]  
 56 represents a Boolean formula as a directed acyclic graph, with terminal nodes labeled by  
 57 Boolean variables and their complements, and with each nonterminal node labeled by a  
 58 Boolean And or Or operation. Restrictions are placed on the structure of the graph such  
 59 that a count of the models can be computed by a single bottom-up traversal. Kimmig et  
 60 al. present a very general *algebraic model counting* [13] framework describing properties of  
 61 Boolean functions that can be efficiently computed from a d-DNNF representation. These  
 62 include standard and weighted model counting, and much more.

63 One shortcoming of existing knowledge compilers is that they have no way to validate that  
 64 their generated results are correct, i.e., that the compiled representation is logically equivalent  
 65 to the original formula. By contrast, all modern SAT solvers can generate checkable proofs  
 66 when they encounter unsatisfiable formulas. The guarantee provided by a checkable certificate  
 67 of correctness enables users of SAT solvers to fully trust their results. Experience has also  
 68 shown that proof production capabilities allow SAT solver developers to quickly detect and  
 69 diagnose bugs in their programs. This, in turn, has led to more reliable SAT solvers.

70 This paper introduces *Partitioned-Operation Graphs* (POGs), a generalization of the  
 71 representations produced by current knowledge compilers. The CPOG (for “certified” POG)  
 72 file format then captures both the structure of a POG and a checkable proof of its logical  
 73 equivalence to a Boolean formula in conjunctive normal form (CNF). A CPOG proof consists  
 74 of a sequence of clause addition and deletion steps, based on an extended resolution proof  
 75 system [19]. We establish a set of conditions that, when satisfied by a CPOG file, guarantee  
 76 that it encodes a well-formed POG and provides a valid equivalence proof.

77 Figure 1 illustrates our certifying knowledge compilation and model counting toolchain.  
 78 Starting with input formula  $\phi_I$ , the D4 knowledge compiler [14] generates a d-DNNF  
 79 representation, and the *proof generator* uses this to generate a CPOG file. The *proof checker*

80 verifies the equivalence of the CNF and CPOG representations. The *ring evaluator* computes  
 81 a standard or weighted model count from the POG representation. As the dashed box in  
 82 Figure 1 indicates, this toolchain moves the root of trust away from the complex and highly  
 83 optimized knowledge compiler to a relatively simple checker and evaluator. Importantly, the  
 84 proof generator need not be trusted—its errors will be caught by the proof checker.

85 To ensure soundness of the abstract CPOG proof system, as well as correctness of its  
 86 concrete implementation, we formally verified the proof system as well as versions of the proof  
 87 checker and ring evaluator in the Lean 4 proof assistant [7]. Running these two programs on  
 88 a particular CPOG file gives assurance that the proof and the model count are correct. Our  
 89 experience with developing a formally verified proof checker has shown that, even within the  
 90 well-understood framework of extended resolution, it can be challenging to formulate a full  
 91 set of requirements that guarantee soundness. In fact, subtle details in the partitioned sum  
 92 rule were identified during this process.

93 We evaluate our toolchain using benchmark formulas from the 2022 standard and weighted  
 94 model competitions. Our tools handle all but the largest d-DNNF graphs generated by D4.  
 95 We also measure the benefits of several optimizations as well as the relative performance of  
 96 the verified checker with one designed for high performance and capacity.

## 97 **2 Related Work**

98 Generating proofs of unsatisfiability in SAT solvers has a long tradition [24] and has become  
 99 widely accepted due to the formulation of clausal proof systems for which proofs can readily  
 100 be generated and efficiently checked [11, 23]. A number of formally verified checkers have  
 101 been developed within different verification frameworks [3, 10, 15, 18]. The associated proofs  
 102 add clauses while preserving satisfiability until the empty clause is derived. By contrast, our  
 103 proofs construct a new propositional formula, and we must verify that it is equivalent to the  
 104 input formula. This requires verifying additional proof steps, including clause deletion steps,  
 105 and subtle invariants, as described in Sections 7 and 10.

106 Fichte, Hecher, and Roland [8] devised the MICE proof framework for model counting  
 107 programs. Their proof rules are based on the algorithms commonly used by model counters.  
 108 They present a series of lemmas and a theorem stating that their checking framework is  
 109 sound. They also present experimental results for cases where they modify existing model  
 110 counters to generate proof traces, and where they generate proof traces from d-DNNF graphs.  
 111 Our toolchain is not directly comparable to theirs—they directly certify the numeric value  
 112 produced by a model counter, while we certify that the graph generated by a knowledge  
 113 compiler is logically equivalent to the input formula. We can then validate a standard or  
 114 weighted count from that representation. Nonetheless, we compare the performance of our  
 115 toolchain to theirs as part of our experimental evaluation in Appendix D.

116 We know of no other work on certifying the output of a knowledge compiler.

## 117 **3 Logical Foundations**

118 Let  $X$  denote a set of Boolean variables, and let  $\alpha$  be an *assignment* of truth values to some  
 119 subset of the variables, where 0 denotes false and 1 denotes true, i.e.,  $\alpha: X' \rightarrow \{0, 1\}$  for  
 120 some  $X' \subseteq X$ . We say the assignment is *total* when it assigns a value to every variable  
 121 ( $X' = X$ ), and that it is *partial* otherwise. The set of all possible total assignments over  $X$  is  
 122 denoted  $\mathcal{U}$ .

123 For each variable  $x \in X$ , we define the *literals*  $x$  and  $\bar{x}$ , where  $\bar{x}$  is the negation of  $x$ .

124 An assignment  $\alpha$  can be viewed as a set of literals, where we write  $\ell \in \alpha$  when  $\ell = x$  and  
 125  $\alpha(x) = 1$  or when  $\ell = \bar{x}$  and  $\alpha(x) = 0$ . We write the negation of literal  $\ell$  as  $\bar{\ell}$ . That is,  $\bar{\bar{\ell}} = \ell$   
 126 when  $\ell = x$  and  $\bar{\bar{\ell}} = x$  when  $\ell = \bar{x}$ .

127 ► **Definition 1** (Boolean Formulas). *The set of Boolean formulas is defined recursively. Each*  
 128 *formula  $\phi$  has an associated dependency set  $\mathcal{D}(\phi) \subseteq X$ , and a set of models  $\mathcal{M}(\phi)$ , consisting*  
 129 *of total assignments that satisfy the formula:*

- 130 1. Boolean constants 0 and 1 are Boolean formulas, with  $\mathcal{D}(0) = \mathcal{D}(1) = \emptyset$ , with  $\mathcal{M}(0) = \emptyset$ ,  
 131 and with  $\mathcal{M}(1) = \mathcal{U}$ .
- 132 2. Variable  $x$  is a Boolean formula, with  $\mathcal{D}(x) = \{x\}$  and  $\mathcal{M}(x) = \{\alpha \in \mathcal{U} \mid \alpha(x) = 1\}$ .
- 133 3. For formula  $\phi$ , its negation, written  $\neg\phi$  is a Boolean formula, with  $\mathcal{D}(\neg\phi) = \mathcal{D}(\phi)$  and  
 134  $\mathcal{M}(\neg\phi) = \mathcal{U} - \mathcal{M}(\phi)$ .
- 135 4. For formulas  $\phi_1, \phi_2, \dots, \phi_k$ , their product  $\phi = \bigwedge_{1 \leq i \leq k} \phi_i$  is a Boolean formula, with  
 136  $\mathcal{D}(\phi) = \bigcup_{1 \leq i \leq k} \mathcal{D}(\phi_i)$  and  $\mathcal{M}(\phi) = \bigcap_{1 \leq i \leq k} \mathcal{M}(\phi_i)$ .
- 137 5. For formulas  $\phi_1, \phi_2, \dots, \phi_k$ , their sum  $\phi = \bigvee_{1 \leq i \leq k} \phi_i$  is a Boolean formula, with  $\mathcal{D}(\phi) =$   
 138  $\bigcup_{1 \leq i \leq k} \mathcal{D}(\phi_i)$  and  $\mathcal{M}(\phi) = \bigcup_{1 \leq i \leq k} \mathcal{M}(\phi_i)$ .

139 We highlight two special classes of Boolean formulas. A formula is in *negation normal*  
 140 *form* when negation is applied only to variables. A formula is in *conjunctive normal form*  
 141 (CNF) when i) it is in negation normal form, and ii) sum is applied only to literals. A CNF  
 142 formula can be represented as a set of *clauses*, each of which is a set of literals. Each clause  
 143 represents the sum of the literals, and the formula is the product of its clauses. We use set  
 144 notation to reference the clauses within a formula and the literals within a clause. A clause  
 145 consisting of a single literal is referred to as a *unit* clause and the literal as a *unit* literal.  
 146 This literal must be assigned value 1 by any satisfying assignment of the formula.

147 ► **Definition 2** (Partitioned-Operation Formula). *A partitioned-operation formula satisfies*  
 148 *the following for all product and sum operations:*

- 149 1. *The arguments to each product must have disjoint dependency sets. That is, operation*  
 150  $\bigwedge_{1 \leq i \leq k} \phi_i$  *requires  $\mathcal{D}(\phi_i) \cap \mathcal{D}(\phi_j) = \emptyset$  for  $1 \leq i < j \leq k$ .*
- 151 2. *The arguments to each sum must have disjoint models. That is, operation  $\bigvee_{1 \leq i \leq k} \phi_i$*   
 152 *requires  $\mathcal{M}(\phi_i) \cap \mathcal{M}(\phi_j) = \emptyset$  for  $1 \leq i < j \leq k$ .*

153 We let  $\wedge^p$  and  $\vee^p$  denote the product and sum operations in a partitioned-operation formula.

## 4 Ring Evaluation of a Boolean Formula

155 We propose a general framework for summarizing properties of Boolean formulas along the  
 156 lines of algebraic model counting [13].

157 ► **Definition 3** (Commutative Ring). *A commutative ring  $\mathcal{R}$  is an algebraic structure*  
 158  $\langle \mathcal{A}, +, \times, \mathbf{0}, \mathbf{1} \rangle$ , *with elements in the set  $\mathcal{A}$  and with commutative and associative opera-*  
 159 *tions  $+$  (addition) and  $\times$  (multiplication), such that multiplication distributes over addition.*  
 160  $\mathbf{0}$  *is the additive identity and  $\mathbf{1}$  is the multiplicative identity. Every element  $a \in \mathcal{A}$  has an*  
 161 *additive inverse  $-a$  such that  $a + -a = \mathbf{0}$ .*

162 We write  $a - b$  as a shorthand for  $a + -b$ .

163 ► **Definition 4** (Ring Evaluation Problem). *For commutative ring  $\mathcal{R}$ , a ring weight function*  
 164 *associates a value  $w(x) \in \mathcal{A}$  with every variable  $x \in X$ . We then define  $w(\bar{x}) \doteq \mathbf{1} - w(x)$ .*

165 *For Boolean formula  $\phi$  and ring weight function  $w$ , the ring evaluation problem computes*

$$166 \quad R(\phi, w) = \sum_{\alpha \in \mathcal{M}(\phi)} \prod_{\ell \in \alpha} w(\ell) \tag{1}$$

167 In this equation, sum  $\Sigma$  is computed using addition operation  $+$ , and product  $\Pi$  is computed  
168 using multiplication operation  $\times$ .

169 Many important properties of Boolean formulas can be expressed as ring evaluation  
170 problems. The (standard) *model counting* problem for formula  $\phi$  requires determining  
171  $|\mathcal{M}(\phi)|$ . It can be cast as a ring evaluation problem by having  $+$  and  $\times$  be addition and  
172 multiplication over rational numbers and using weight function  $w(x) = 1/2$  for every variable  
173  $x$ . Ring evaluation of formula  $\phi$  gives the *density* of the formula, i.e., the fraction of all  
174 possible total assignments that are models. For  $n = |X|$ , scaling the density by  $2^n$  yields the  
175 number of models. This formulation avoids the need for a “smoothing” operation, in which  
176 redundant expressions are inserted into the formula [6].

177 The *weighted model counting* problem is also defined over rational numbers. Some  
178 formulations allow independently assigning weights  $W(x)$  and  $W(\bar{x})$  for each variable  $x$   
179 and its complement, with the possibility that  $W(x) + W(\bar{x}) \neq 1$ . We can cast this as a  
180 ring evaluation problem by letting  $r(x) = W(x) + W(\bar{x})$ , performing ring evaluation with  
181 weight function  $w(x) = W(x)/r(x)$  for each variable  $x$ , and computing the weighted count  
182 as  $\mathbf{R}(\phi, w) \times \prod_{x \in X} r(x)$ . Of course, this requires that  $r(x) \neq 0$  for all  $x \in X$ .

183 The *function hashing problem* provides a test of inequivalence for Boolean formulas.  
184 That is, for  $n = |X|$ , let  $\mathcal{R}$  be a finite field with  $|\mathcal{A}| = m$  such that  $m \geq 2n$ . For each  
185  $x \in X$ , chose a value from  $\mathcal{A}$  at random for  $w(x)$ . Two formulas  $\phi_1$  and  $\phi_2$  will clearly have  
186  $\mathbf{R}(\phi_1, w) = \mathbf{R}(\phi_2, w)$  if they are logically equivalent. If they are not equivalent, then the  
187 probability that  $\mathbf{R}(\phi_1, w) \neq \mathbf{R}(\phi_2, w)$  will be at least  $(1 - \frac{1}{m})^n \geq (1 - \frac{1}{2n})^n > 1/2$ . Function  
188 hashing can therefore be used as part of a randomized algorithm for equivalence testing [2].  
189 For example, it can compare different runs on a single formula, either from different compilers  
190 or from a single compiler with different configuration parameters.

## 191 5 Partitioned-Operation Graphs (POGs)

192 Performing ring evaluation on an arbitrary Boolean formula could be intractable, but it is  
193 straightforward for a formula with partitioned operations:

194 ► **Proposition 5.** *Ring evaluation with operations  $\neg$ ,  $\wedge^P$ , and  $\vee^P$  satisfies the following for  
195 any weight function  $w$ :*

$$196 \quad \mathbf{R}(\neg\phi, w) = \mathbf{1} - \mathbf{R}(\phi, w) \quad (2)$$

$$197 \quad \mathbf{R}\left(\bigwedge_{1 \leq i \leq k}^P \phi_i, w\right) = \prod_{1 \leq i \leq k} \mathbf{R}(\phi_i, w) \quad (3)$$

$$198 \quad \mathbf{R}\left(\bigvee_{1 \leq i \leq k}^P \phi_i, w\right) = \sum_{1 \leq i \leq k} \mathbf{R}(\phi_i, w) \quad (4)$$

199 As is described in Section 10, we have proved these three equations using Lean 4.

200 A *partitioned-operation graph* (POG) is a directed, acyclic graph with nodes  $N$  and edges  
201  $E \subseteq N \times N$ . We denote nodes with boldface symbols, such as  $\mathbf{u}$  and  $\mathbf{v}$ . When  $(\mathbf{u}, \mathbf{v}) \in E$ ,  
202 node  $\mathbf{v}$  is said to be a *child* of node  $\mathbf{u}$ . The in- and out-degrees of node  $\mathbf{u}$  are defined  
203 as  $\text{indegree}(\mathbf{u}) = |E \cap (N \times \{\mathbf{u}\})|$ , and  $\text{outdegree}(\mathbf{u}) = |E \cap (\{\mathbf{u}\} \times N)|$ . Node  $\mathbf{u}$  is said  
204 to be *terminal* if  $\text{outdegree}(\mathbf{u}) = 0$ . A terminal node is labeled by a Boolean constant or  
205 variable. Node  $\mathbf{u}$  is said to be *nonterminal* if  $\text{outdegree}(\mathbf{u}) > 0$ . A nonterminal node is  
206 labeled by one of the three Boolean operations:  $\neg$ ,  $\wedge^P$  or  $\vee^P$ . Node  $\mathbf{u}$  can be labeled with  $\neg$   
207 only if  $\text{outdegree}(\mathbf{u}) = 1$ . It can be labeled with operation  $\wedge^P$  or  $\vee^P$  only if it satisfies the  
208 partitioning restriction for that operation. Every POG has a designated *root node*  $\mathbf{r}$ .

209 In essence, a POG represents a partitioned-operation formula with a sharing of common  
 210 subformulas. Indeed, every node in the graph can be viewed as a partitioned-operation  
 211 formula, and so we write  $\phi_{\mathbf{u}}$  as the formula denoted by node  $\mathbf{u}$ . Each such formula has a set  
 212 of models, and we write  $\mathcal{M}(\mathbf{u})$  as a shorthand for  $\mathcal{M}(\phi_{\mathbf{u}})$ .

213 We define the *size* of POG  $P$ , written  $|P|$ , to be the the number of nodes labeled  $\wedge^P$  or  
 214  $\vee^P$  plus the number of edges from these nodes to their children. Ring evaluation of  $P$  can be  
 215 performed with at most  $|P|$  ring operations by traversing the graph from the terminal nodes  
 216 up to the root, computing a value  $\mathbf{R}(\mathbf{u}, w)$  for each node  $\mathbf{u}$ . The final result is then  $\mathbf{R}(\mathbf{r}, w)$ .

217 POGs generalize the d-DNNF graphs in two ways:

- 218 ■ They allow negation of arbitrary nodes in the graph, not just variables.
  - 219 ■ They allow arbitrary arguments to a  $\vee^P$  operation, as long as these have disjoint models.
- 220 By contrast, each sum node  $\mathbf{u}$  in a d-DNNF must have two children  $\mathbf{u}_1$  and  $\mathbf{u}_0$ , and for  
 221 these there must be a *decision variable*  $x$  such that any total assignment  $\alpha \in \mathcal{M}(\mathbf{u}_b)$  has  
 222  $\alpha(x) = b$ , for  $b \in \{0, 1\}$ .

223 These generalizations allow more flexibility in the POG representation while maintaining the  
 224 ability to efficiently perform ring evaluation.

## 225 6 Clausal Proof Framework

226 We write (possibly subscripted)  $\theta$  for formulas encoded as clauses, possibly with extension  
 227 variables. We write (possibly subscripted)  $\phi$  for formulas that use no extension variables.

228 A proof in our framework consists of a sequence of clause addition and deletion steps,  
 229 with each step preserving the set of solutions to the original formula. The status of the proof  
 230 at any step is represented as a set of *active* clauses  $\theta$ , i.e., those that have been added but  
 231 not yet deleted. Our framework is based on *extended* resolution [19], where proof steps can  
 232 introduce new *extension variables* encoding Boolean formulas over input and prior extension  
 233 variables. Let  $Z$  denote the set of extension variables occurring in formula  $\theta$ . Starting with  $\theta$   
 234 equal to input formula  $\phi_I$ , the proof must maintain the invariant that  $\phi_I \Leftrightarrow \exists Z \theta$ .

235 Clauses can be added in two different ways. One is when they serve as the *defining clauses*  
 236 for an extension variable. This form occurs only when defining  $\wedge^P$  and  $\vee^P$  operations, as is  
 237 described in Section 7. Clauses can also be added or deleted based on *implication redundancy*.  
 238 That is, when clause  $C$  satisfies  $\theta \Rightarrow C$  for formula  $\theta$ , then it can either be added to  $\theta$  to  
 239 create the formula  $\theta \cup \{C\}$  or it can be deleted from  $\theta \cup \{C\}$  to create  $\theta$ .

240 We use *reverse unit propagation* (RUP) to certify implication redundancy when adding or  
 241 deleting clauses [9, 22]. RUP is the core rule supported by standard proof checkers [11, 23] for  
 242 propositional logic. It provides a simple and efficient way to check a sequence of applications  
 243 of the resolution proof rule [17]. Let  $C = \{\ell_1, \ell_2, \dots, \ell_p\}$  be a clause to be proved redundant  
 244 with respect to formula  $\theta$ . Let  $D_1, D_2, \dots, D_k$  be a sequence of supporting *antecedent* clauses,  
 245 such that each  $D_i$  is in  $\theta$ . A RUP step proves that  $\bigwedge_{1 \leq i \leq k} D_i \Rightarrow C$  by showing that the  
 246 combination of the antecedents plus the negation of  $C$  leads to a contradiction. The negation  
 247 of  $C$  is the formula  $\bar{\ell}_1 \wedge \bar{\ell}_2 \wedge \dots \wedge \bar{\ell}_p$ , having a CNF representation consisting of  $p$  unit clauses  
 248 of the form  $\bar{\ell}_i$  for  $1 \leq i \leq p$ . A RUP check processes the clauses of the antecedent in sequence,  
 249 inferring additional unit clauses. In processing clause  $D_i$ , if all but one of the literals in the  
 250 clause is the negation of one of the accumulated unit clauses, then we can add this literal to  
 251 the accumulated set. That is, all but this literal have been falsified, and so it must be set to  
 252 true for the clause to be satisfied. The final step with clause  $D_k$  must cause a contradiction,  
 253 i.e., all of its literals are falsified by the accumulated unit clauses.

254 Compared to the proofs of unsatisfiability generated by SAT solvers, ours have important

■ **Table 1** CPOG Step Types.  $C$ : clause identifier,  $L$ : literal,  $V$ : variable

		Rule		Description
$C$	<b>a</b>	$L^* 0$	$C^+ 0$	Add RUP clause
	<b>d</b>	$C$	$C^+ 0$	Delete RUP clause
$C$	<b>p</b>	$V L^* 0$		Declare $\wedge^P$ operation
$C$	<b>s</b>	$V L L$	$C^+ 0$	Declare $\vee^P$ operation
	<b>r</b>	$L$		Declare root literal

255 differences. Most significantly, each proof step must preserve the set of solutions with  
 256 respect to the input variables; our proofs must therefore justify both clause deletions and  
 257 additions. By contrast, an unsatisfiability proof need only guarantee that no proof step  
 258 causes a satisfiable set of clauses to become unsatisfiable, and therefore it need only justify  
 259 clause additions.

## 260 **7 The CPOG POG Representation and Proof System**

261 A CPOG file provides both a declaration of a POG, as well as a checkable proof that a  
 262 Boolean formula, given in conjunctive normal form, is logically equivalent to the POG.  
 263 The proof format draws its inspiration from the LRAT [10] and QRAT [12] formats for  
 264 unquantified and quantified Boolean formulas, respectively. Key properties include:

- 265 ■ The file contains declarations of  $\wedge^P$  and  $\vee^P$  operations to describe the POG. Declaring a  
 266 node  $\mathbf{u}$  implicitly adds an *extension* variable  $u$  and a set of *defining* clauses  $\theta_u$  encoding  
 267 the product or sum operation.
- 268 ■ Boolean negation is supported implicitly by allowing the arguments of the  $\vee^P$  and  $\wedge^P$   
 269 operations to be literals and not just variables.
- 270 ■ The file contains explicit clause addition steps. A clause can only be added if it is logically  
 271 implied by the existing clauses. A sequence of clause identifiers must be listed as a *hint*  
 272 providing a RUP verification of the implication.
- 273 ■ The file contains explicit clause deletion steps. A clause can only be deleted if it is  
 274 logically implied by the remaining clauses. A sequence of clause identifiers must be listed  
 275 as a *hint* providing a RUP verification of the implication.
- 276 ■ The checker must track the dependency set for every input and extension variable. For  
 277 each  $\wedge^P$  operation, the checker must ensure that the dependency sets for its arguments  
 278 are disjoint. The associated extension variable has a dependency set equal to the union  
 279 of those of its arguments.
- 280 ■ Declaring a  $\vee^P$  operation requires a sequence of clauses providing a RUP proof that  
 281 the arguments are mutually exclusive. Only binary  $\vee^P$  operations are allowed to avoid  
 282 requiring multiple proofs of disjointness

### 283 **7.1 Syntax**

284 Table 1 shows the declarations that can occur in a CPOG file. As with other clausal proof  
 285 formats, a variable is represented by a positive integer  $v$ , with the first ones being input  
 286 variables and successive ones being extension variables. Literal  $\ell$  is represented by a signed  
 287 integer, with  $-v$  being the logical negation of variable  $v$ . Each clause is indicated by a positive  
 288 integer identifier  $C$ , with the first ones being the IDs of the input clauses and successive ones

■ **Table 2** Defining Clauses for Product (A) and Sum (B) Operations

(A). Product Operation $\wedge^P$						(B). Sum Operation $\vee^P$			
ID	Clause					ID	Clause		
$i$	$v$	$-\ell_1$	$-\ell_2$	$\dots$	$-\ell_k$	$i$	$-v$	$\ell_1$	$\ell_2$
$i+1$	$-v$	$\ell_1$				$i+1$	$v$	$-\ell_1$	
$i+2$	$-v$	$\ell_2$				$i+2$	$v$	$-\ell_2$	
		$\dots$							
$i+k$	$-v$	$\ell_k$							

289 being the IDs of added clauses. Clause identifiers must be totally ordered, such that any  
 290 clause identifier  $C'$  given in the hint when adding clause  $C$  must have  $C' < C$ .

291 The first set of proof rules are similar to those in other clausal proofs. Clauses can be  
 292 added via RUP addition (command **a**), with a sequence of antecedent clauses (the “hint”).  
 293 Similarly for clause deletion (command **d**).

294 The declaration of a *product* operation, creating a node with operation  $\wedge^P$ , has the form:

295 
$$i \quad \mathbf{p} \quad v \quad \ell_1 \quad \ell_2 \quad \dots \quad \ell_k \quad 0$$

296 Integer  $i$  is a new clause ID,  $v$  is a positive integer that does not correspond to any previous  
 297 variable, and  $\ell_1, \ell_2, \dots, \ell_k$  is a sequence of  $k$  integers representing literals of existing variables.  
 298 As Table 2(A) shows, this declaration implicitly causes  $k + 1$  clauses to be added to the  
 299 proof, defining extension variable  $v$  to be the product of its arguments.

300 The dependency sets for the arguments represented by each pair of literals  $\ell_i$  and  $\ell_j$  must  
 301 be disjoint, for  $1 \leq i < j \leq k$ . A product operation may have no arguments, representing  
 302 Boolean constant 1. The only clause added to the proof will be the unit literal  $v$ . A reference  
 303 to literal  $-v$  then provides a way of representing constant 0.

304 The declaration of a *sum* operation, creating a node with operation  $\vee^P$ , has the form:

305 
$$i \quad \mathbf{s} \quad v \quad \ell_1 \quad \ell_2 \quad H \quad 0$$

306 Integer  $i$  is a new clause ID,  $v$  is a positive integer that does not correspond to any previous  
 307 variable, and  $\ell_1$  and  $\ell_2$  are signed integers representing literals of existing variables. Hint  $H$   
 308 consists of a sequence of clause IDs, all of which must be defining clauses for other POG  
 309 operations. As Table 2(B) shows, this declaration implicitly causes three clauses to be added  
 310 to the proof, defining extension variable  $v$  to be the sum of its arguments. The hint must  
 311 provide a RUP proof of the clause  $\bar{\ell}_1 \vee \bar{\ell}_2$ , showing that the two children of this node have  
 312 disjoint models.

313 Finally, the literal denoting the root of the POG is declared with the **r** command. It  
 314 can occur anywhere in the file. Except in degenerate cases, it will be the extension variable  
 315 representing the root of a graph.

## 316 7.2 Semantics

317 The defining clauses for a product or sum operation uniquely define the value of its extension  
 318 variable for any assignment of values to the argument variables. For the extension variable  $u$   
 319 associated with any POG node  $\mathbf{u}$ , we can therefore prove that any total assignment  $\alpha$  to the  
 320 input variables which also satisfies the POG definitions will assign a value to  $u$  such that  
 321  $\alpha(u) = 1$  if and only if  $\alpha \in \mathcal{M}(\mathbf{u})$ .



322 The sequence of operator declarations, asserted clauses, and clause deletions represents  
 323 a systematic transformation of the input formula into a POG formula. Validating all of  
 324 these steps serves to prove that the POG is logically equivalent to the input formula. At the  
 325 completion of the proof, the following FINAL CONDITIONS must hold:

- 326 1. There is exactly one remaining clause that was added via RUP addition, and this is a  
 327 unit clause consisting of root literal  $r$ .
- 328 2. All of the input clauses have been deleted.

329 In other words, at the end of the proof it must hold that the active clauses be exactly those  
 330 in  $\theta_P \doteq \{\{r\}\} \cup \bigcup_{u \in P} \theta_u$ , the formula consisting of unit clause  $\{r\}$  and the defining clauses  
 331 of the POG nodes. Recognizing that any total assignment to the input variables implicitly  
 332 defines the assignments to the extension variables, we can see that  $\theta_P$  is the clausal encoding  
 333 of  $\phi_r$ . Let  $\phi_I$  denote the input formula. The sequence of clause addition steps provides a  
 334 *forward implication* proof that  $\mathcal{M}(\phi_I) \subseteq \mathcal{M}(\phi_r)$ . That is, any total assignment  $\alpha$  satisfying  
 335 the input formula must also satisfy the formula represented by the POG. Conversely, each  
 336 proof step that deletes an input clause proves that any total assignment  $\alpha$  that falsifies the  
 337 clause must falsify  $\phi_r$ . Deleting all but the final asserted clause and all input clauses provides  
 338 a *reverse implication* proof that  $\mathcal{M}(\phi_r) \subseteq \mathcal{M}(\phi_I)$ .

339 Appendix A shows the CPOG description for an input formula with five clauses, yielding  
 340 a POG with five nonterminal nodes. It explains how the clause addition and deletion steps  
 341 yield a proof of equivalence between the input formula and its POG representation.

## 342 **8** Generating CPOG from d-DNNF

343 A d-DNNF graph can be directly translated into a POG, although in doing this conversion, our  
 344 program also performs simplifications to eliminate Boolean constants. Except in degenerate  
 345 cases, we can therefore assume that the POG does not contain any constant nodes. In  
 346 addition, negation is only applied to variables, and so we can merge the negation operations  
 347 into the terminal nodes, viewing the POG as consisting of *literal* nodes corresponding to  
 348 input variables and their negations, along with *nonterminal* nodes, which can be further  
 349 classified as *product* and *sum* nodes.

### 350 **8.1 Forward Implication Proof**

351 For input formula  $\phi_I$  and its translation into a POG  $P$  with root node  $\mathbf{r}$ , the most challenging  
 352 part of the proof is to show that  $\mathcal{M}(\phi_I) \subseteq \mathcal{M}(\phi_r)$ , i.e., that any total assignment  $\alpha$  that is  
 353 a model of  $\phi_I$  and the POG definition clauses will yield  $\alpha(r) = 1$ , for root literal  $r$ . This  
 354 part of the proof consists of a series of clause assertions leading to one adding  $\{r\}$  as a unit  
 355 clause. We have devised two methods of generating this proof. The *monolithic* approach  
 356 makes just one call to a proof-generating SAT solver and has it determine the relationship  
 357 between the two representations. The *structural* approach recursively traverses the POG,  
 358 generating proof obligations at each node encountered. It may require multiple calls to a  
 359 proof-generating SAT solver.

360 As notation, let  $\psi$  be a subset of the clauses in  $\phi_I$ . For partial assignment  $\rho$ , the expression  
 361  $\psi|_\rho$  denotes the set of clauses  $\gamma$  obtained from  $\psi$  by: i) eliminating any clause containing a  
 362 literal  $\ell$  such that  $\rho(\ell) = 1$ , ii) for the remaining clauses eliminating those literals  $\ell$  for which  
 363  $\rho(\ell) = 0$ , and iii) eliminating any duplicate clauses. In doing these simplifications, we also  
 364 track the *provenance* of each simplified clause  $C$ , i.e., which of the (possibly multiple) input  
 365 clauses simplified to become  $C$ . More formally, for  $C \in \psi|_\rho$ , we let  $\text{Prov}_\rho(C, \psi)$  denote those

366 clauses  $C' \in \psi$ , such that  $C' \subseteq C \cup \bigcup_{\ell \in \rho} \bar{\ell}$ . We then extend the definition of  $\text{Prov}$  to any  
 367 simplified formula  $\gamma$  as  $\text{Prov}_\rho(\gamma, \psi) = \bigcup_{C \in \gamma} \text{Prov}_\rho(C, \psi)$ .

368 The monolithic approach takes advantage of the clausal representations of the input  
 369 formula  $\phi_I$  and the POG formula  $\phi_r$ . We can express the negation of  $\phi_r$  in clausal form  
 370 as  $\theta_{\bar{r}} \doteq \bigcup_{\mathbf{u} \in P} \theta_{\mathbf{u}}|_{\{\bar{r}\}}$ . Forward implication will hold when  $\phi_I \Rightarrow \phi_r$ , or equivalently when  
 371 the formula  $\phi_I \wedge \theta_{\bar{r}}$  is unsatisfiable, where the conjunction can be expressed as the union  
 372 of the two sets of clauses. The proof generator writes the clauses to a file and invokes a  
 373 proof-generating SAT solver. For each clause  $C$  in the unsatisfiability proof, it adds clause  
 374  $\{r\} \cup C$  to the CPOG proof, and so the empty clause in the proof becomes the unit clause  
 375  $\{r\}$ . Our experimental results show that this approach can be very effective and generates  
 376 short proofs for smaller problems, but it does not scale well enough for routine use.

377 We describe the structural approach to proof generation as a recursive procedure  
 378  $\text{validate}(\mathbf{u}, \rho, \psi)$  taking as arguments POG node  $\mathbf{u}$ , partial assignment  $\rho$ , and a set of clauses  
 379  $\psi \subseteq \phi_I$ . The procedure adds a number of clauses to the proof, culminating with the addition  
 380 of the *target* clause:  $u \vee \bigvee_{\ell \in \rho} \bar{\ell}$ , indicating that  $(\bigwedge_{\ell \in \rho} \ell) \Rightarrow u$ , i.e., that any total assignment  
 381  $\alpha$  such that  $\rho \subseteq \alpha$  will assign  $\alpha(u) = 1$ . The top-level call has  $\mathbf{u} = \mathbf{r}$ ,  $\rho = \emptyset$ , and  $\psi = \phi_I$ .  
 382 The result will therefore be to add unit clause  $\{r\}$  to the proof. Here we present a correct,  
 383 but somewhat inefficient formulation of  $\text{validate}$ . We then refine it with some optimizations.

384 The recursive call  $\text{validate}(u, \rho, \psi)$  assumes that we have traversed a path from the root  
 385 node down to node  $\mathbf{u}$ , with the literals encountered in the product nodes forming the partial  
 386 assignment  $\rho$ . The set of clauses  $\psi$  can be a proper subset of the input clauses  $\phi_I$  when a  
 387 product node has caused a splitting into clauses containing disjoint variables. The subgraph  
 388 with root node  $\mathbf{u}$  should be a POG representation of the formula  $\psi|_\rho$ .

389 The process for generating such a proof depends on the form of node  $\mathbf{u}$ :

- 390 **1.** If  $\mathbf{u}$  is a literal  $\ell'$ , then the formula  $\psi|_\rho$  must consist of the single unit clause  $C = \{\ell'\}$ ,  
 391 such that any  $C' \in \text{Prov}_\rho(C, \psi)$  must have  $C' \subseteq \{\ell'\} \cup \bigcup_{\ell \in \rho} \bar{\ell}$ . Any of these can serve as  
 392 the target clause.
- 393 **2.** If  $\mathbf{u}$  is a sum node with children  $\mathbf{u}_1$  and  $\mathbf{u}_0$ , then, since the node originated from a d-DNNF  
 394 graph, there must be some variable  $x$  such that either  $\mathbf{u}_1$  is a literal node for  $x$  or  $\mathbf{u}_1$  is a  
 395 product node containing a literal node for  $x$  as a child. In either case, we recursively call  
 396  $\text{validate}(\mathbf{u}_1, \rho \cup \{x\}, \psi)$ . This will cause the addition of the target clause  $u_1 \vee \bar{x} \vee \bigvee_{\ell \in \rho} \bar{\ell}$ .  
 397 Similarly, either  $\mathbf{u}_0$  is a literal node for  $\bar{x}$  or  $\mathbf{u}_0$  is a product node containing a literal  
 398 node for  $\bar{x}$  as a child. In either case, we recursively call  $\text{validate}(\mathbf{u}_0, \rho \cup \{\bar{x}\}, \psi)$ , causing  
 399 the addition of the target clause  $u_0 \vee x \vee \bigvee_{\ell \in \rho} \bar{\ell}$ . These recursive results can be combined  
 400 with the second and third defining clauses for  $\mathbf{u}$  (see Table 2(B)) to generate the target  
 401 clause for  $\mathbf{u}$ .
- 402 **3.** If  $\mathbf{u}$  is a product node, then we can divide its children into a set of literal nodes  $L$  and a  
 403 set of nonterminal nodes  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$ .
  - 404 **a.** For each literal  $\ell \in L$ , we must prove that any total assignment  $\alpha$ , such that  $\rho \subset \alpha$  has  
 405  $\alpha(\ell) = 1$ . In some cases, this can be done by simple Boolean constraint propagation  
 406 (BCP). In other cases, we must prove that the formula  $\psi|_{\rho \cup \{\bar{\ell}\}}$  is unsatisfiable. We do  
 407 so by writing the formula to a file, invoking a proof-generating SAT solver, and then  
 408 converting the generated unsatisfiability proof into a sequence of clause additions in  
 409 the CPOG file.
  - 410 **b.** For a single nonterminal child ( $k = 1$ ), we recursively call  $\text{validate}(\mathbf{u}_1, \rho \cup \bigcup_{\ell \in L} \ell, \psi)$ .
  - 411 **c.** For multiple nonterminal children ( $k > 1$ ), it must be the case that the clauses in  
 412  $\gamma = \psi|_\rho$  can be partitioned into  $k$  subsets  $\gamma_1, \gamma_2, \dots, \gamma_k$  such that  $\mathcal{D}(\gamma_i) \cap \mathcal{D}(\gamma_j) = \emptyset$  for  
 413  $1 \leq i < j \leq k$ , and we can match each node  $\mathbf{u}_i$  to subset  $\gamma_i$  based on its literals. For each

414  $i$  such that  $1 \leq i \leq k$ , let  $\psi_i = \text{Prov}_\rho(\gamma_i, \psi)$ , i.e., those input clauses in  $\psi$  that, when  
 415 simplified, became clause partition  $\gamma_i$ . We recursively call `validate` ( $\mathbf{u}_i, \rho \cup \bigcup_{\ell \in L} \ell, \psi_i$ ).  
 416 We then generate the target clause for node  $\mathbf{u}$ , creating the hint by combining the results  
 417 from the BCP and SAT calls for the literals, the recursively computed target clauses, and  
 418 all but the first defining clause for node  $\mathbf{u}$  (see Table 2(A)).

419 Most of these steps involve a polynomial number of operations per recursive call, with the  
 420 exception of those that call a SAT solver to validate a literal.

## 421 8.2 Reverse Implication Proof

422 Completing the equivalence proof of input formula  $\phi_I$  and its POG representation with root  
 423 node  $\mathbf{r}$  requires showing that  $\mathcal{M}(\phi_{\mathbf{r}}) \subseteq \mathcal{M}(\phi_I)$ . This is done in the CPOG framework by  
 424 first deleting all asserted clauses, except for the final unit clause for root literal  $r$ , and then  
 425 deleting all of the input clauses.

426 The asserted clauses can be deleted in reverse order, using the same hints that were used  
 427 in their original assertions. By reversing the order, those clauses that were used in the hint  
 428 when a clause was added will still remain when it is deleted.

429 Each input clause deletion can be done as a single RUP step. The proof generator  
 430 constructs the hint sequence from the defining clauses of the POG nodes via a single, bottom-  
 431 up pass through the graph. The RUP deletion proof for input clause  $C$  effectively proves  
 432 that any total assignment  $\alpha$  that satisfies the POG definition clauses but does not satisfy  
 433  $C$  will yield  $\alpha(r) = 0$ . It starts with the set of literals  $\{\bar{\ell} \mid \ell \in C\}$ , describing the required  
 434 condition for assignment  $\alpha$  to falsify clause  $C$ . It then adds literals via unit propagation  
 435 until a conflict arises. Unit literal  $r$  gets added right away, setting up a potential conflict.

436 Working upward through the graph, node  $\mathbf{u}$  is *marked* when the collected set of literals  
 437 force  $u$  to evaluate to 0. When marking  $\mathbf{u}$ , the program adds  $\bar{u}$  to the RUP literals and adds  
 438 the appropriate defining clause to the hint. A literal node for  $\ell$  will be marked if  $\ell \in C$ , with  
 439 no hint required. If product node  $\mathbf{u}$  has some child  $\mathbf{u}_i$  that is marked, then  $\mathbf{u}$  is marked and  
 440 clause  $i + 1$  from among its defining clauses (see Table 2(A)) is added to the hint. Marking  
 441 sum node  $\mathbf{u}$  requires that its two children are marked. The first defining clause for this node  
 442 (see Table 2(B)) will then be added to the hint. At the very end, the program (assuming the  
 443 reverse implication holds) will attempt to mark root node  $\mathbf{r}$ , which would require  $\alpha(r) = 0$ ,  
 444 yielding a conflict.

445 It can be seen that the reverse implication proof will be polynomial in the size of the  
 446 POG, because each clause deletion requires a single RUP step having a hint with length  
 447 bounded by the number of POG nodes.

## 448 9 Optimizations

449 The performance of the structural proof generator for forward implication, both in its  
 450 execution time and the size of the proof generated, can be improved by two optimizations  
 451 described here. A key feature is that they do not require any changes to the proof framework—  
 452 they build on the power of extended resolution to enable new logical structures. They involve  
 453 declaring new product nodes to encode products of literals. These nodes are not part of the  
 454 POG representation of the formula; they serve only to enable the forward implication proof.  
 455 Here we summarize the two optimizations. Appendix B provides more details.

456 *Literal Grouping:* A single recursive step of `validate` can encounter product nodes having  
 457 many—tens or even hundreds—of literals as children. The earlier formulation of `validate`

458 considers each literal  $\ell \in L$  separately, calling a SAT solver for every literal that cannot  
 459 be validated with BCP. Literal grouping handles all of these literals together. It defines  
 460 product node  $\mathbf{v}$  having the literals as children. The goal then becomes to prove that any  
 461 total assignment must yield 1 for extension variable  $v$ . Calling a solver with  $v$  set to 0 yields  
 462 an unsatisfiability proof that can be mapped back to a sequence of clause additions in the  
 463 CPOG file validating all of the literals.

464 *Lemmas:* Our formulation of `validate` requires each call at a node  $\mathbf{u}$  to recursively validate  
 465 all of its children. This effectively expands the graph into a tree, potentially requiring an  
 466 exponential number of recursive steps. Instead, for each node  $\mathbf{u}$  having  $\text{indegree}(\mathbf{u}) > 1$ , the  
 467 program can define and generate the proof of a *lemma* for  $\mathbf{u}$  when it is first reached by a call  
 468 to `validate` and then apply this lemma for this and subsequent calls. The lemma states that  
 469 the POG with root node  $\mathbf{u}$  satisfies forward implication for a formula  $\gamma_{\mathbf{u}}$ , where some of the  
 470 clauses in this formula are input clauses from  $\phi_I$ , but others are simplified versions of input  
 471 clauses. The key idea is to introduce product nodes to encode (via DeMorgan’s Laws) the  
 472 simplified clauses and have these serve as lemma arguments.

473 The combination of these two optimization guarantees that i) each call to `validate` for  
 474 a product node will cause at most one invocation of the SAT solver, and ii) each call to  
 475 `validate` for any node  $\mathbf{u}$  will cause further recursive calls only once. Our experimental results  
 476 (Appendix D) shows that these optimizations yield substantial benefits.

## 477 10 A Formally Verified Toolchain

478 We set out to formally verify the system with two goals in mind: first, to ensure that the  
 479 CPOG framework is mathematically sound; and second, to implement correct-by-construction  
 480 proof checking and ring evaluation (the “Trusted Code” components of Figure 1). These two  
 481 goals are achieved with a single proof development in the Lean 4 programming language [7].  
 482 Verification was greatly aided by the Aesop [16] automated tactic. In this section, we briefly  
 483 describe the functionality we implemented and what we proved about it. More information  
 484 is provided in Appendix C.

485 **Proof checking.** The goal of a CPOG proof is to construct a POG that is equivalent to  
 486 the input CNF  $\phi_I$ . The checker begins by parsing the input formula, initializing the set of  
 487 active clauses to  $\theta \leftarrow \phi_I$ , and initializing the POG  $P$  to an empty one. It then processes  
 488 every step of the CPOG proof, either modifying its state by adding/deleting clauses in  $\theta$  and  
 489 adding nodes to  $P$ , or throwing an exception if a step is incorrect. Afterwards, it carries  
 490 out the FINAL CONDITIONS check of Section 7.2. Throughout the process, we maintain  
 491 invariants which ensure that  $P$  is partitioned and that a successful final check entails the  
 492 logical equivalence of  $\phi_I$  and  $\phi_{\mathbf{r}}$ , where  $\mathbf{r}$  is the final POG root.

493 The specifications we use to state these invariants are built on a general theory of  
 494 propositional logic, mirroring Section 3. Following the DIMACS CNF convention, we define  
 495 the data types `Var` of variables being positive natural numbers, `ILit` of literals being non-zero  
 496 integers, and `PropForm Var` of propositional formulas. Assignments of truth values are taken to  
 497 be total functions `PropAssignment Var := Var → Bool`. Requiring totality is not a limitation:  
 498 instead of talking about two equal, partial assignments to a subset  $X' \subseteq X$  of variables, we  
 499 can more conveniently talk about two total assignments that agree on  $X'$ . We write  $\sigma \models \varphi$   
 500 when  $\sigma : \text{PropAssignment Var}$  satisfies  $\varphi : \text{PropForm Var}$ .

501 The invariants refer to the checker state `st` with fields `st.inputCnf` for  $\phi_I$ , `st.clauseDb` for  
 502  $\theta$ , `st.pog` for  $P$ , `st.pogDefsForm` for the POG definitions formula  $\bigwedge_{\mathbf{u} \in P} \theta_{\mathbf{u}}$ , and `st.allVars`  
 503 for all variables (original and extension) added so far. For any  $\mathbf{u} \in P$ , `st.pog.toPropForm u`

504 computes  $\phi_{\mathbf{u}}$ . The first two invariants state that assignments to original variables extend  
 505 uniquely to extension variables defining the POG nodes. In the formalization, we split this  
 506 into extension and uniqueness:

```

507 /-- Any assignment satisfying  $\varphi_1$  extends to  $\varphi_2$  while preserving values on X. -/
508 def extendsOver (X : Set Var) ( $\varphi_1 \varphi_2$  : PropForm Var) :=
509    $\forall (\sigma_1 : \text{PropAssignment Var}), \sigma_1 \models \varphi_1 \rightarrow \exists \sigma_2, \sigma_1.\text{agreeOn } X \sigma_2 \wedge \sigma_2 \models \varphi_2$ 
510 /-- Assignments satisfying  $\varphi$  are determined on Y by their values on X. -/
511 def uniqueExt (X Y : Set Var) ( $\varphi$  : PropForm Var) :=
512    $\forall (\sigma_1 \sigma_2 : \text{PropAssignment Var}), \sigma_1 \models \varphi \rightarrow \sigma_2 \models \varphi \rightarrow \sigma_1.\text{agreeOn } X \sigma_2 \rightarrow$ 
513      $\sigma_1.\text{agreeOn } Y \sigma_2$ 
514
515
516 invariants.extends_pogDefsForm : extendsOver st.inputCnf.vars  $\top$  st.pogDefsForm
517 invariants.uep_pogDefsForm : uniqueExt st.inputCnf.vars st.allVars st.pogDefsForm

```

519 The next one guarantees that the set of solutions over the original variables is preserved:

```

520
521 def equivalentOver (X : Set Var) ( $\varphi_1 \varphi_2$  : PropForm Var) :=
522   extendsOver X  $\varphi_1 \varphi_2 \wedge \text{extendsOver } X \varphi_2 \varphi_1$ 
523
524 invariants.equivInput : equivalentOver st.inputCnf.vars st.inputCnf st.clauseDb
525

```

526 Finally, for every node  $\mathbf{u} \in P$  with corresponding literal  $u$  we ensure that  $\phi_{\mathbf{u}}$  is partitioned  
 527 (Definition 2) and relate  $\phi_{\mathbf{u}}$  to its clausal encoding  $\theta_u \doteq u \wedge \bigwedge_{v \in P} \theta_v$ :

```

528
529 def partitioned : PropForm Var  $\rightarrow$  Prop
530 | tr | fls | var _ => True
531 | neg  $\varphi$  =>  $\varphi.\text{partitioned}$ 
532 | disj  $\varphi \psi$  =>  $\varphi.\text{partitioned} \wedge \psi.\text{partitioned} \wedge \forall \tau, \neg(\tau \models \varphi \wedge \tau \models \psi)$ 
533 | conj  $\varphi \psi$  =>  $\varphi.\text{partitioned} \wedge \psi.\text{partitioned} \wedge \varphi.\text{vars} \cap \psi.\text{vars} = \emptyset$ 
534
535 invariants.partitioned :  $\forall u : \text{ILit}, (\text{st.pog.toPropForm } u).\text{partitioned}$ 
536 invariants.equivalent_lits :  $\forall u : \text{ILit}, \text{equivalentOver } \text{st.inputCnf.vars}$ 
537    $(u \wedge \text{st.pogDefsForm}) (\text{st.pog.toPropForm } x)$ 
538

```

539 These are maintained by valid CPOG proofs. Together with a few additional invariants  
 540 that ensure the correctness of cached computations, they imply the soundness theorem for  
 541  $P$  with root node  $\mathbf{r}$ : the FINAL CONDITIONS ensure that  $\theta$  is logically equivalent to  $\theta_P \doteq$   
 542  $\{\{r\}\} \cup \bigcup_{\mathbf{u} \in P} \theta_u$  (equivalent st.clauseDb ( $\mathbf{r} \wedge \text{st.pogDefsForm}$ )), and from the invariants  
 543 we can prove  $\phi_I \iff \phi_{\mathbf{r}}$  (equivalent st.inputCnf (st.pog.toPropForm  $\mathbf{r}$ )). After certifying  
 544 a proof, the checker can pass its in-memory POG representation to the ring evaluator.

545 **Ring evaluation.** We formalized the central quantity (1) in the ring evaluation problem  
 546 (Definition 4) in a commutative ring  $R$  as follows:

```

547
548 def weightSum {R : Type} [CommRing R]
549   (weight : Var  $\rightarrow$  R) ( $\varphi$  : PropForm Var) (s : Finset Var) : R :=
550    $\sum \tau \text{ in models } \varphi \text{ s}, \prod x \text{ in } s, \text{if } \tau x \text{ then weight } x \text{ else } 1 - \text{weight } x$ 
551

```

552 The rules for efficient ring evaluation of partitioned formulas are expressed as:

```

553
554 def ringEval (weight : Var  $\rightarrow$  R) : PropForm Var  $\rightarrow$  R
555 | tr => 1
556 | fls => 0
557 | var x => weight x
558 | neg  $\varphi$  => 1 - ringEval weight  $\varphi$ 

```

## XX:14 Certified Knowledge Compilation

```
559 | disj  $\varphi$   $\psi$  => ringEval weight  $\varphi$  + ringEval weight  $\psi$ 
560 | conj  $\varphi$   $\psi$  => ringEval weight  $\varphi$  * ringEval weight  $\psi$ 
```

562 Proposition 5 is then formalized as follows:

```
563
564 theorem ringEval_eq_weightSum (weight : Var → R) { $\varphi$  : PropForm Var} :
565   partitioned  $\varphi$  → ringEval weight  $\varphi$  = weightSum weight  $\varphi$  (vars  $\varphi$ )
```

567 To efficiently compute the ring evaluation of a formula represented by a POG node, we  
568 implemented `Pog.ringEval` and then proved that it matches the specification above:

```
569
570 theorem ringEval_eq_ringEval (pog : Pog) (weight : Var → R) (x : Var) :
571   pog.ringEval weight x = (pog.toPropForm x).ringEval weight
```

573 Applying this to the output of our verified CPOG proof checker, which we know to be  
574 partitioned and equivalent to the input formula  $\phi_I$ , we obtain a proof that our toolchain  
575 computes the correct ring evaluation of  $\phi_I$ .

576 **Model counting.** Finally, we established that ring evaluation with the appropriate weights  
577 corresponds to the standard model count. To do so, we defined a function that carries out  
578 an integer calculation of the number of models over a set of variables of cardinality `nVars`:

```
579
580 def countModels (nVars : Nat) : PropForm Var → Nat
581 | tr      => 2^nVars
582 | fls    => 0
583 | var _   => 2^(nVars - 1)
584 | neg  $\varphi$  => 2^nVars - countModels nVars  $\varphi$ 
585 | disj  $\varphi$   $\psi$  => countModels nVars  $\varphi$  + countModels nVars  $\psi$ 
586 | conj  $\varphi$   $\psi$  => countModels nVars  $\varphi$  * countModels nVars  $\psi$  / 2^nVars
```

588 We then formally proved that for a partitioned formula whose variables are among a finite  
589 set `s`, this computation really does count the number of models over `s`:

```
590
591 theorem countModels_eq_card_models { $\varphi$  : PropForm Var} {s : Finset Var} :
592   vars  $\varphi$   $\subseteq$  s → partitioned  $\varphi$  → countModels (card s)  $\varphi$  = card (models  $\varphi$  s)
```

594 In particular, taking `s` to be exactly the variables of  $\varphi$ , we have that the number of models  
595 of  $\varphi$  on its variables is `countModels  $\varphi$  (card (vars  $\varphi$ ))`.

## 596 11 Implementations

597 We have implemented versions of the three programs that, along with the D4 knowledge  
598 compiler, form the toolchain illustrated in Figure 1<sup>1</sup>. The proof generator is the same in  
599 both cases, since it need not be trusted. Our *verified* versions of the proof checker and ring  
600 evaluator have been formally verified within the Lean 4 theorem prover. Our long term  
601 goal is to rely on these. Our *prototype* versions are written in C (checker) and Python (ring  
602 evaluator). These are currently faster and more scalable, but we anticipate their need will  
603 diminish as the verified version is further optimized.

604 Our proof generator is written in C/C++ and uses CADICAL [1] as its SAT solver. To  
605 convert proof steps back into hinted CPOG clause additions, the generator can use either its  
606 own RUP proof generator, or it can invoke DRAT-TRIM [11]. The latter yields shorter proofs

---

<sup>1</sup> All tools, along with a set of benchmark problems are available at <https://github.com/rebryant/cpog>

607 and scales well to large proofs, but each invocation has a high startup cost. We therefore  
608 only use it when solving larger problems (currently ones with over 1000 clauses).

609 The proof generator can optionally be instructed to generate a *one-sided* proof, providing  
610 only the reverse-implication portion of the proof via input clause deletion. This can provide  
611 useful information—any assignment that is a model for the compiled representation must  
612 also be a model for the input formula—even when full validation is impractical.

613 Our prototype ring evaluator can perform both standard and weighted model counting.  
614 It performs arithmetic over a subset of the rationals we call  $\mathbf{Q}_{2,5}$ , consisting of numbers of  
615 the form  $a \cdot 2^b \cdot 5^c$ , for integers  $a$ ,  $b$ , and  $c$ , and with  $a$  implemented to have arbitrary range.  
616 Allowing scaling by powers of 2 enables the density computation and rescaling required for  
617 standard model counting. Allowing scaling by powers of both 2 and 5 enables exact decimal  
618 arithmetic, handling the weights used in the weighted model counting competitions. To give  
619 a sense of scale, the counter generated a result with 260,909 decimal digits for one of the  
620 weighted benchmarks.

## 621 12 Experimental Evaluation

622 We summarize the results of our experiments here. Appendix D provides a more complete  
623 description. For our evaluation, we used the public benchmark problems from the 2022  
624 standard and weighted model competitions. We found that there were 180 unique CNF files  
625 among these, ranging in size from 250 to 2,753,207 clauses. With a runtime limit of 4000  
626 seconds, D4 completed for 124 of the benchmark problems. Our proof generator was able to  
627 convert all of these into POGs, with their declarations ranging from 304 to 2,761,457,765  
628 (median 848,784) defining clauses. We ran our proof generator with a time limit of 10,000  
629 seconds. It was able to generate full proofs for 108 of the problems and one-sided proofs  
630 for an additional 9 of them, leaving just 7 with no verification. The prototype checker  
631 successfully verified all of the generated proofs. The longest runtime for the combination of  
632 proof generator and checker for a full proof was 13,145 seconds. It is worth noting that, to  
633 date, we have not found any errors in the d-DNNF graphs generated by D4.

634 We found that the monolithic approach for generating the forward implication proof  
635 works well for smaller POGs, but it becomes unreliable and inefficient for larger ones. We  
636 also found that our two optimizations: literal grouping and lemmas can provide substantial  
637 improvements in proof size and runtime. Running the verified proof checker in Lean 4 showed  
638 the same scaling with respect to the proof size as did the prototype checker, but it required,  
639 on average, around 5.9 times longer.

640 Finally, the runtimes for the MICE toolchain versus our CPOG toolchain showed little  
641 correlation, reflecting the fact that the two solve different problems and use different ap-  
642 proaches. In general, our CPOG toolchain showed better scaling, in part due to its ability to  
643 control the recursion through lemmas.

## 644 13 Conclusions

645 We are hopeful that having checkable proofs for knowledge compilers will allow them to  
646 be used in applications where high levels of trust are required, and that it will provide a  
647 useful tool for developers of knowledge compilers. Our experiments demonstrate that our  
648 toolchain can already handle problems nearly at the limits of current knowledge compilers.  
649 Further engineering and optimization of our proof generator and checker could improve their  
650 performance and capacity substantially.

651 ——— **References** ———

- 652 1 Armin Biere. CaDiCaL at the SAT Race 2019. In *Proc. of SAT Race 2019 – Solver and*  
653 *Benchmark Descriptions*, volume B-2019-1 of *Department of Computer Science Series of*  
654 *Publications B*, pages 8–9. University of Helsinki, 2019.
- 655 2 Manuel Blum, Ashok K. Chandra, and Mark N. Wegman. Equivalence of free Boolean  
656 graphs can be decided probabilistically in polynomial time. *Information Processing Letters*,  
657 10(2):80–82, 18 March 1980.
- 658 3 Luís Cruz-Filipe, Marijn J. H. Heule, Warren A. Hunt, Matt Kaufmann, and Peter Schneider-  
659 Kamp. Efficient certified RAT verification. In *Conference on Automated Deduction (CADE)*,  
660 volume 10395 of *LNCS*, pages 220–236, 2017.
- 661 4 Adnan Darwiche. A compiler for deterministic, decomposable negation normal form. In  
662 *Association for the Advancement of Artificial Intelligence (AAAI)*, 2002.
- 663 5 Adnan Darwiche. New advances in compiling CNF to decomposable negation normal form. In  
664 *European Conference on Artificial Intelligence*, pages 328–332, 2004.
- 665 6 Adnan Darwiche and Pierre Marquis. A knowledge compilation map. *Journal of Artificial*  
666 *Intelligence Research*, 17, 2002.
- 667 7 Leonardo de Moura and Sebastian Ulrich. The Lean 4 theorem prover and programming  
668 language. In *Conference on Automated Deduction (CADE)*, volume 12699 of *LNAI*, pages  
669 625–635, 2021.
- 670 8 Johannes K Fichte, Markus Hecher, and Valentin Roland. Proofs for propositional model  
671 counting. In *Theory and Applications of Satisfiability Testing (SAT)*. Schloss Dagstuhl-Leibniz-  
672 Zentrum für Informatik, 2022.
- 673 9 Evgueni I. Goldberg and Yakov Novikov. Verification of proofs of unsatisfiability for CNF  
674 formulas. In *Design, Automation and Test in Europe (DATE)*, pages 886–891, 2003.
- 675 10 Marijn J. H. Heule, Warren A. Hunt, Matt Kaufmann, and Nathan D. Wetzler. Efficient,  
676 verified checking of propositional proofs. In *Interactive Theorem Proving*, volume 10499 of  
677 *LNCS*, pages 269–284, 2017.
- 678 11 Marijn J. H. Heule, Warren A. Hunt, Jr., and Nathan D. Wetzler. Verifying refutations with  
679 extended resolution. In *Conference on Automated Deduction (CADE)*, volume 7898 of *LNCS*,  
680 pages 345–359, 2013.
- 681 12 Marijn J. H. Heule, Martina Seidl, and Armin Biere. A unified proof system for QBF  
682 preprocessing. In *International Joint Conference on Automated Reasoning (IJCAR)*, volume  
683 8562 of *LNCS*, pages 91–106, 2014.
- 684 13 Angelika Kimmig, Guy Van den Broeck, and Luc De Raedt. Algebraic model counting. *Journal*  
685 *of Applied Logic*, 22:46–62, July 2017.
- 686 14 Jean-Marie Lagniez and Pierre Marquis. An improved decision-DNNF compiler. In *Interna-*  
687 *tional Joint Conference on Artificial Intelligence*, pages 667–673, 2017.
- 688 15 Peter Lammich. Efficient verified (UN)SAT certificate checking. *J. Autom. Reason.*, 64(3):513–  
689 532, 2020.
- 690 16 Jannis Limperg and Asta Halkjær From. Aesop: White-box best-first proof search for Lean.  
691 In *Certified Programs and Proofs (CPP)*, pages 253–266. ACM, 2023.
- 692 17 John A. Robinson. A machine-oriented logic based on the resolution principle. *J.ACM*,  
693 12(1):23–41, January 1965.
- 694 18 Yong Kiam Tan, Marijn J. H. Heule, and Magnus O. Myreen. cake\_lpr: Verified propagation  
695 redundancy checking in CakeML. In *Tools and Algorithms for the Construction and Analysis*  
696 *of Systems (TACAS), Part II*, volume 12652 of *LNCS*, pages 223–241, 2021.
- 697 19 G. S. Tseitin. On the complexity of derivation in propositional calculus. In *Automation of*  
698 *Reasoning: 2: Classical Papers on Computational Logic 1967–1970*, pages 466–483. Springer,  
699 1983.
- 700 20 Sebastian Ullrich and Leonardo de Moura. Counting immutable beans: reference counting  
701 optimized for purely functional programming. In *Implementation and Application of Functional*  
702 *Languages (IFL)*, pages 3:1–3:12. ACM, 2019.



- 703 21 Leslie G. Valiant. The complexity of enumeration and reliability problems. *SIAM Journal of*  
704 *Computing*, 8(3):410–421, 1979.
- 705 22 Allen Van Gelder. Verifying RUP proofs of propositional unsatisfiability. In *Proc. of the 10th*  
706 *Int. Symposium on Artificial Intelligence and Mathematics (ISAIM 2008)*, 2008.
- 707 23 Nathan D. Wetzler, Marijn J. H. Heule, and Warren A. Hunt Jr. DRAT-trim: Efficient checking  
708 and trimming using expressive clausal proofs. In *Theory and Applications of Satisfiability*  
709 *Testing (SAT)*, volume 8561 of *LNCS*, pages 422–429, 2014.
- 710 24 Lintao Zhang and Sharad Malik. Validating SAT solvers using an independent resolution-based  
711 checker: Practical implementations and other applications. In *Design, Automation and Test*  
712 *in Europe (DATE)*, pages 880–885, 2003.

713 **A CPOG Example**

714 Figure 2 illustrates an example formula and shows how the CPOG file declares its POG  
 715 representation. The input formula (A) consists of five clauses over variables  $x_1$ ,  $x_2$ ,  $x_3$ , and  
 716  $x_4$ . The generated POG (B) has six nonterminal nodes representing four products and two  
 717 sums. We name these by the node type (product **p** or sum **s**), subscripted by the ID of the  
 718 extension variable. The first part of the CPOG file (C) declares these nodes using clause IDs  
 719 that increment by three or four, depending on whether the node has two children or three.  
 720 The last two nonzero values in the sum declarations are the hint providing the required  
 721 mutual exclusion proof.

722 **A.1 Node Declarations**

723 We step through portions of the file to provide a better understanding of the CPOG proof  
 724 framework. Figure 2(D) shows the defining clauses that are implicitly defined by the POG  
 725 operation declarations. These do not appear in the CPOG file. Referring back to the  
 726 declarations of the sum nodes in Figure 2(C), we can see that the declaration of node **s**<sub>7</sub> had  
 727 clause IDs 7 and 10 as the hint. We can see in Figure 2(A) that these two clauses form a  
 728 RUP proof for the clause  $\bar{p}_5 \vee \bar{p}_6$ , showing that the two children of **s**<sub>7</sub> have disjoint models.  
 729 Similarly, node **s**<sub>10</sub> is declared as having clause IDs 16 and 19 as the hint. These form a RUP  
 730 proof for the clause  $\bar{p}_8 \vee \bar{p}_9$ , showing that the two children of **s**<sub>10</sub> have disjoint models.

731 **A.2 Forward Implication Proof**

732 Figure 2(E) provides the sequence of assertions leading to unit clause 36, consisting of the  
 733 literal  $s_{10}$ . This clause indicates that **s**<sub>10</sub> is implied by the input clauses, i.e., any total  
 734 assignment  $\alpha$  satisfying the input clauses must have  $\alpha(s_{10}) = 1$ . Working backward, we  
 735 can see that clause 29 indicates that variable  $p_8$  will be implied by the input clauses when  
 736  $\alpha(x_1) = 0$ , while clause 34 indicates that node  $p_9$  will be implied by the input clauses when  
 737  $\alpha(x_1) = 1$ . These serve as the hint for clause 36.

738 **A.3 Reverse Implication Proof**

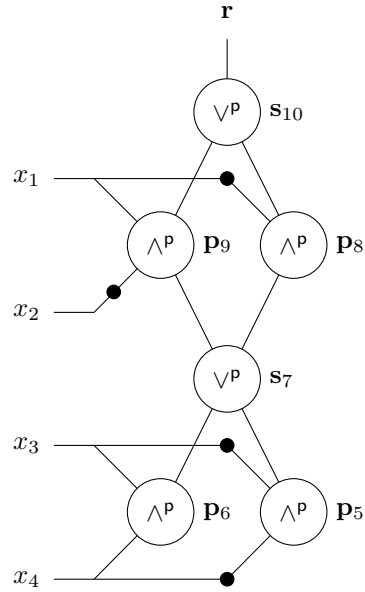
739 Figure 2(F) shows the RUP proof steps required to delete the input clauses. Consider the first  
 740 of these, deleting input clause  $\bar{x}_1 \vee x_3 \vee \bar{x}_4$ . The requirement is to show that there is no total  
 741 assignment  $\alpha$  that falsifies this clause but assigns  $\alpha(s_{10}) = 1$ . The proof proceeds by first  
 742 assuming that the clause is false, requiring  $\alpha(x_1) = 1$ ,  $\alpha(x_3) = 0$ , and  $\alpha(x_4) = 1$ . The hint  
 743 then consist of unit clauses (e.g., clause 36 asserting that  $\alpha(p_{10}) = 1$ ) or clauses that cause  
 744 unit propagation. Hint clauses 8 and 10 force the assignments  $\alpha(p_5) = \alpha(p_6) = 0$ . These,  
 745 plus hint clause 12 force  $\alpha(s_7) = 0$ . This, plus hint clauses 16 and 21 force  $\alpha(p_8) = \alpha(p_9) = 0$ ,  
 746 leading via clause 22 to  $\alpha(s_{10}) = 0$ . But this contradicts clause 36, completing the RUP  
 747 proof. The deletion hints for the other input clauses follow similar patterns—they work from  
 748 the bottom nodes of the POG upward, showing that any total assignment that falsifies the  
 749 clause must assign  $\alpha(s_{10}) = 0$ .

750 Deleting the asserted clauses is so simple that we do not show it. It involves simply  
 751 deleting the clauses from clause number 35 down to clause number 25, with each deletion  
 752 using the same hint as were used to add that clause. In the end, therefore, only the defining  
 753 clauses for the POG nodes and the unit clause asserting  $s_{10}$  remain, completing a proof that  
 754 the POG is logically equivalent to the input formula.

(A) Input Formula

ID	Clauses
1	-1 3 -4 0
2	-1 -3 4 0
3	3 -4 0
4	1 -3 4 0
5	-1 -2 0

(B) POG Representation



(C) POG Declaration

ID	CPOG line	Explanation
6	p 5 -3 -4 0	$p_5 = \bar{x}_3 \wedge^P \bar{x}_4$
9	p 6 3 4 0	$p_6 = x_3 \wedge^P x_4$
12	s 7 5 6 7 10 0	$s_7 = p_5 \vee^P p_6$
15	p 8 -1 7 0	$p_8 = \bar{x}_1 \wedge^P s_7$
18	p 9 1 -2 7 0	$p_9 = x_1 \wedge^P \bar{x}_2 \wedge^P s_7$
22	s 10 8 9 16 19 0	$s_{10} = p_8 \vee^P p_9$
	r 10	Root $r = s_{10}$

(D) Implicit Clauses

ID	Clauses	Explanation
6	5 3 4 0	Define $p_5$
7	-5 -3 0	
8	-5 -4 0	
9	6 -3 -4 0	Define $p_6$
10	-6 3 0	
11	-6 4 0	
12	-7 5 6 0	Define $s_7$
13	7 -5 0	
14	7 -6 0	
15	8 1 -7 0	Define $p_8$
16	-8 -1 0	
17	-8 7 0	
18	9 -1 2 -7 0	Define $p_9$
19	-9 1 0	
20	-9 -2 0	
21	-9 7 0	
22	-10 8 9 0	Define $s_{10}$
23	10 -8 0	
24	10 -9 0	

(E) CPOG Assertions

ID	Clause	Hint	Explanation
25	a 5 1 3 0	3 6 0	$\bar{x}_1 \wedge \bar{x}_3 \Rightarrow p_5$
26	a 6 1 -3 0	4 9 0	$\bar{x}_1 \wedge x_3 \Rightarrow p_6$
27	a 3 7 1 0	13 25 0	$\bar{x}_3 \wedge \bar{x}_1 \Rightarrow s_7$
28	a 7 1 0	27 14 26 0	$\bar{x}_1 \Rightarrow s_7$
29	a 8 1 0	28 15 0	$\bar{x}_1 \Rightarrow p_8$
30	a 5 -1 3 0	1 6 0	$x_1 \wedge \bar{x}_3 \Rightarrow p_5$
31	a 6 -1 -3 0	2 9 0	$x_1 \wedge x_3 \Rightarrow p_6$
32	a 3 7 -1 0	13 30 0	$\bar{x}_3 \wedge x_1 \Rightarrow s_7$
33	a 7 -1 0	32 14 31 0	$x_1 \Rightarrow s_7$
34	a 9 -1 0	5 33 18 0	$x_1 \Rightarrow p_9$
35	a 1 10 0	23 29 0	$\bar{x}_1 \Rightarrow s_{10}$
36	a 10 0	35 24 34 0	$s_{10}$

(F) Input Clause Deletions

CPOG line	Explanation
d 1 36 8 10 12 16 21 22 0	Delete clause 1
d 2 36 7 11 12 16 21 22 0	Delete clause 2
d 3 36 8 10 12 17 19 22 0	Delete clause 3
d 4 36 7 11 12 17 19 22 0	Delete clause 4
d 5 36 16 20 22 0	Delete clause 5

■ **Figure 2** Example formula (A), its POG representation (B), and its CPOG proof (C), (E), and (F)

755 Observe that the forward implication proof shown in Figure 2 must “visit” nodes  $\mathbf{p}_5$  and  
 756  $\mathbf{p}_6$  twice, separately considering assignments where  $\alpha(x_1) = 1$  (clauses 25 and 26) and where  
 757  $\alpha(x_1) = 1$  (clauses 30 and 31). Section B.3 shows how to use a lemma such that these nodes  
 758 are each only visited once.

## 759 **B Optimizations for Forward Implication Proofs**

760 The two optimizations we have implemented exploit the power of our general resolution  
 761 framework to define new structures within the proof. They create new product nodes that  
 762 are not part of the POG representation.

### 763 **B.1 Literal Grouping**

764 A single recursive step of `validate` can encounter product nodes having many literals as  
 765 children. The naive formulation of `validate` considers each literal  $\ell \in L$  separately. Literal  
 766 grouping allows all literals to be validated with a single call to a SAT solver. It collects  
 767 those literals  $\ell_1, \ell_2, \dots, \ell_m$  that cannot be validated by BCP and defines a product node  $\mathbf{v}$   
 768 having these literals as children. The goal then becomes to prove that any total assignment  
 769 must yield 1 for extension variable  $v$ . A single call to the solver can generate this proof by  
 770 invoking it on the formula  $\psi|_{\rho} \cup \theta_v|_{\{\bar{v}\}}$ , which should be unsatisfiable. The proof steps can  
 771 be mapped back into clause addition steps in the CPOG file, incorporating the input clauses  
 772 and the defining clauses for  $\mathbf{v}$  into the hints.

### 773 **B.2 Lemmas**

774 As we have noted, the recursive calling of `validate` starting at root  $\mathbf{r}$  effectively expands  
 775 the POG into a tree, and this can lead to an exponential number of calls. These shared  
 776 subgraphs arise when the knowledge compiler employs *clause caching* to detect that the  
 777 simplified set of clauses arising from one partial assignment to the literals matches that of a  
 778 previous partial assignment [4]. When this d-DNNF node is translated into POG node  $\mathbf{u}$ ,  
 779 the proof generator can assume (and also check), that there is a simplified set of clauses  $\gamma_{\mathbf{u}}$   
 780 for which the subgraph with root  $\mathbf{u}$  is its POG representation.

781 The proof generator can exploit the sharing of subgraphs by constructing and proving a  
 782 *lemma* for each node  $\mathbf{u}$  having  $\text{indegree}(\mathbf{u}) > 1$ . This proof shows that any total assignment  
 783  $\alpha$  that satisfies formula  $\gamma_{\mathbf{u}}$  and the defining clauses for the POG must yield  $\alpha(u) = 1$ . This  
 784 lemma is then invoked for every node having  $\mathbf{u}$  as a child. As a result, the generator will  
 785 make recursive calls during a call to `validate` only once for each node in the POG.

786 The challenge for implementing this strategy is to find a way to represent the clauses  
 787 for the simplified formula  $\gamma_{\mathbf{u}}$  in the CPOG file. Some may be unaltered input clauses, and  
 788 these can be used directly. Others, however will be clauses that do not appear in the input  
 789 formula. We implement these by adding POG product nodes to the CPOG file to create the  
 790 appropriate clauses. Consider an *argument* clause  $C \in \gamma_{\mathbf{u}}$  with  $C = \ell_1 \vee \ell_2 \vee \dots \vee \ell_k$ . If we  
 791 define product node  $\mathbf{v}$  with arguments  $\bar{\ell}_1, \bar{\ell}_2, \dots, \bar{\ell}_k$ , then by Table 2(A), we will introduce  
 792 a defining clause  $v \vee \ell_1 \vee \ell_2 \vee \dots \vee \ell_k$ . We call this a *synthetic* clause having  $\bar{v}$  as the *guard*  
 793 *literal*. That is, a partial assignment  $\rho$  such that  $\rho(v) = 0$  will *activate* the clause, causing it  
 794 to represent argument clause  $C$ . On the other hand, a partial assignment with  $\rho(v) = 1$  will  
 795 cause the clause to become a tautology and therefore have no effect.

796 Suppose for every clause  $C_j \in \gamma_{\mathbf{u}}$  that does not correspond to an input clause, we generate  
 797 a synthetic clause  $C'_j$  with guard literal  $\bar{v}_j$ , for  $1 \leq j \leq m$ . Let  $\gamma'_{\mathbf{u}}$  be the formula where each

798 clause  $C_j$  is replaced by synthetic clause  $C'_j$ , while input clauses in  $\gamma_{\mathbf{u}}$  are left unchanged.  
 799 Let  $\beta = \{\bar{v}_1, \bar{v}_2, \dots, \bar{v}_m\}$ . Invoking `validate`( $\mathbf{u}, \beta, \gamma'_{\mathbf{u}}$ ) will then prove a lemma, given by the  
 800 target clause  $u \vee v_1 \vee v_2 \vee \dots \vee v_m$ , showing that any total assignment  $\alpha$  that activates the  
 801 synthetic clauses will have  $\alpha(u) = 1$ .

802 Later, when node  $\mathbf{u}$  is encountered by a call to `validate`( $\mathbf{u}, \rho, \psi$ ), we invoke the lemma  
 803 by showing that each synthetic clause  $C_j$  matches some simplified clause in  $\psi|_{\rho}$ . More  
 804 precisely, for  $1 \leq j \leq m$ , we use clause addition to assert the clause  $\bar{v}_j \vee \bigvee_{\ell \in \rho} \bar{\ell}$ , showing that  
 805 synthetic clause  $C_j$  will be activated. Combining the lemma with these activations provides  
 806 a derivation of the target clause for the call to `validate`.

807 Observe that the lemma structure can be hierarchical, since a shared subgraph may  
 808 contain nodes that are themselves roots of shared subgraphs. Nonetheless, the principles  
 809 described allow the definition, proof, and applications of a lemma for each shared node in  
 810 the graph. For any node  $\mathbf{u}$ , the first call to `validate`( $\mathbf{u}, \rho, \psi$ ) may require further recursion,  
 811 but any subsequent call can simply reuse the lemma proved by the first call.

### 812 B.3 Lemma Example

813 Figure 3 shows an alternate forward implication proof for the example of Figure 2 using  
 814 a lemma to represent the shared node  $s_7$ . We can see that the POG with this node as  
 815 root encodes the Boolean formula  $x_3 \leftrightarrow x_4$ , having a CNF representation consisting of the  
 816 clauses  $\{x_3, \bar{x}_4\}$  and  $\{\bar{x}_3, x_4\}$ . The product node declarations shown in Figure 3(A) create  
 817 synthetic clauses 25 and 28 to encode these arguments with activating literals  $\bar{v}_{11}$  and  $\bar{v}_{12}$ ,  
 818 respectively. Clauses 31–34 then provide a proof of the lemma, stating that any assignment  
 819  $\alpha$  that activates these clauses will assign 1 to  $s_7$ . Clauses 35 and 36 state that an assignment  
 820 with  $\alpha(x_1) = 0$  will cause the first synthetic clause to activate due to input clause 3, and it  
 821 will cause the second synthetic clause to activate due to input clause 4. From this, clause 37  
 822 can use the lemma to state that assigning 0 to  $x_1$  will cause  $s_7$  to evaluate to 1. Similarly,  
 823 clauses 39 and 40 serve to activate the synthetic clauses when  $\alpha(x_1) = 1$ , due to input clauses  
 824 1 and 2, and clause 41 then uses the lemma to state that assigning 1 to  $x_1$  will cause  $s_7$   
 825 to evaluate to 1.

826 In this example, adding the lemma increases the proof length, but that is only because it  
 827 is such a simple formula.

## 828 C More Information about the Formally Verified Checker

829 Our verified toolchain has been implemented in Lean 4, which is based on a formal logical  
 830 foundation in which expressions have a computational interpretation. As with other proof  
 831 assistants like Isabelle and Coq, a function defined within the formal system can be compiled  
 832 to efficient code. At the same time, we can state and prove claims about the function within  
 833 the system, thereby verifying that the functions compute the intended results.

834 Our code is contained in the directory `ProofChecker`. The main checker loop is imple-  
 835 mented in `Checker/CheckerCore.lean`. That file defines a structure `PreState` that contains  
 836 all the relevant data structures, which include the input CNF formula, the clause database,  
 837 the POG under construction, and the root literal. We define a `State` of the checker to be  
 838 a `PreState` that satisfies all the invariants that we need to establish the final result, and  
 839 the bulk of our work involved showing that these are indeed preserved by each step of the  
 840 proof, which modifies the clause database and the POG. The ring evaluator and model  
 841 counter are contained in the `Count` directory. Here we provide a few additional notes on the  
 842 implementation and clarify the sense in which the toolchain has been verified.

(A) Additional nodes

ID	CPOG line	Explanation
25	p 11 -3 4 0	$v_{11} = \bar{x}_3 \wedge^p x_4$
28	p 12 3 -4 0	$v_{12} = x_3 \wedge^p \bar{x}_4$

(B) Implicit Clauses

ID	Clauses	Explanation
25	11 3 -4 0	Argument clause $\{x_3, \bar{x}_4\}$ , activated by $\bar{v}_{11}$
26	-11 -3 0	
27	-11 4 0	
28	12 -3 4 0	Argument clause $\{\bar{x}_3, x_4\}$ , activated by $\bar{v}_{12}$
29	-12 3 0	
30	-12 -4 0	

(C) CPOG Assertions

ID	Clause	Hint	Explanation
Lemma Proof			
31	a 5 11 12 3 0	25 6 0	$(\bar{v}_{11} \wedge \bar{v}_{12}) \wedge \bar{x}_3 \Rightarrow p_5$
32	a 6 11 12 -3 0	28 9 0	$(\bar{v}_{11} \wedge \bar{v}_{12}) \wedge x_3 \Rightarrow p_6$
33	a 3 7 11 12 0	13 31 0	$(\bar{v}_{11} \wedge \bar{v}_{12}) \wedge \bar{x}_3 \Rightarrow s_7$
34	a 7 11 12 0	33 14 32 0	$(\bar{v}_{11} \wedge \bar{v}_{12}) \Rightarrow s_7$
Lemma Application #1			
35	a -11 1 0	26 27 3 0	$\bar{x}_1 \Rightarrow \bar{v}_{11}$
36	a -12 1 0	29 30 4 0	$\bar{x}_1 \Rightarrow \bar{v}_{12}$
37	a 7 1 0	35 36 34 0	$\bar{x}_1 \Rightarrow s_7$
38	a 8 1 0	37 15 0	$\bar{x}_1 \Rightarrow p_8$
Lemma Application #2			
39	a -11 -1 0	26 27 1 0	$x_1 \Rightarrow \bar{v}_{11}$
40	a -12 -1 0	29 30 2 0	$x_1 \Rightarrow \bar{v}_{12}$
41	a 7 -1 0	39 40 34 0	$x_1 \Rightarrow s_7$
42	a 9 -1 0	5 41 18 0	$x_1 \Rightarrow p_9$
43	a 1 10 0	23 38 0	$\bar{x}_1 \Rightarrow s_{10}$
44	a 10 0	43 24 42 0	$s_{10}$

■ **Figure 3** Example of lemma definition, proof, and application

## 843 C.1 Implementation notes

844 When thinking about the formal verification, it is helpful to distinguish between the data  
 845 structures that play a role in the code that is executed and the definitions that serve as  
 846 a mathematical reference, allowing us to state our specifications and prove that they are  
 847 met. Among the former is our representation of CNF formulas: a literal is represented as a  
 848 nonzero integer, a clause is an array of literals, and a CNF formula is an array of clauses.

```
849 def ILit := { i : Int // i ≠ 0 }
850 abbrev IClause := Array ILit
851 abbrev ICnf := Array IClause
```

854 We define the clause database `ClauseDb` as a hashmap that stores each clause together with  
 855 a flag indicating whether it has been deleted. Each element of a POG is either a variable, a  
 856 binary disjunction (sum), or an arbitrary conjunction (product):

```
857 inductive PogElt where
858   | var : Var → PogElt
859   | disj : Var → ILit → ILit → PogElt
860   | conj : Var → Array ILit → PogElt
```

863 In the first case, the argument `x` in the expression `var x` is the index of an original variable;  
 864 in `disj x left right` and `conj x args` it is an extension variable appearing in the CPOG  
 865 file. Note that representing edges as literals allows us to negate the arguments to `disj` and  
 866 `conj`. A `Pog` is then an array of `PogElt` that is well-founded in the sense that each element  
 867 depends only on prior elements in the array.

868 On the mathematical side, we define propositional formulas and their semantics in the usual  
 869 way. Functions `ILit.toPropForm`, `IClause.toPropForm`, `ICnf.toPropForm`, and `Pog.toPropForm`  
 870 relate our data structures to propositional formulas and their semantics. For example, given  
 871 a literal `u`, `Pog.toPropForm P u` denotes the (unnegated or negated) interpretation of the  
 872 node corresponding to `u` in the POG `P` as a propositional formula  $\phi_u / \neg\phi_u$  over the original  
 873 variables. Lean provides us with convenient “anonymous projector” notation that allows us  
 874 to write `P.toPropForm u` instead of `Pog.toPropForm P u` when `P` has type `Pog`, `C.toPropForm`  
 875 instead of `IClause.toPropForm C` when `C` has type `IClause`, etc.

876 In order to reason about the behavior of the checker, we found it important to work  
 877 with propositional formulas modulo logical equivalence, a structure known in logic as the  
 878 *Lindenbaum–Tarski algebra*. We defined the quotient and lifted the Boolean operations and  
 879 the entailment relation to this new type. The advantage is that equivalent formulas give rise  
 880 to equal elements in the quotient. This makes it much easier to prove equivalences involving  
 881 complicated expressions, since it enables us to substitute elements of the quotient for others  
 882 even when the corresponding formulas are merely equivalent but not equal. This gave rise to  
 883 auxiliary challenges, however. For example, it is no longer straightforward to say that an  
 884 element of the quotient “depends” on a variable, since equivalent formulas can have different  
 885 variables—consider  $x \vee \neg x$  and  $\top$ . Instead, we made use of a semantic notion of dependence,  
 886 in which an element  $\phi$  of the quotient depends on a variable if and only if for some truth  
 887 assignment the value of  $\phi$  changes when the assignment to that variable is flipped.

## 888 C.2 Trust

889 In this subsection, we clarify what has been verified and what has to be trusted. Recall that  
 890 our first step is to parse CNF and CPOG files in order to read in the initial formula and the  
 891 proof. We do not verify this step. Instead, the verified checker exposes flags `--print-cnf`

892 and `--print-cpog` which reprint the consumed formula or proof, respectively. Comparing  
 893 this to the actual files using `diff` provides an easy way of ensuring that what was parsed  
 894 matches their contents. This involves trusting only the correctness of the print procedure  
 895 and `diff`. Similarly, if one wants to establish the correctness of the POG contained in the  
 896 CPOG file, one can print out the POG that is constructed by the checker and compare.

897 Lean’s code extraction replaces calculations on natural numbers and integers with efficient  
 898 but unverified arbitrary precision versions. Lean also uses an efficient implementation of  
 899 arrays; within the formal system, these are defined in terms of lists, but code extraction  
 900 replaces them with dynamic arrays and uses reference counting to allow destructive updates  
 901 when it is safe to do so [20]. Finally, Lean’s standard library implements hashmaps in terms  
 902 of arrays. Many of the basic properties of hashmaps have been formally verified, but not all.  
 903 In particular, we make use of a fold operation whose verification is not yet complete. Thus  
 904 our proofs depend on the assumption that the fold operation has the expected properties.

905 In summary, in addition to trusting Lean’s foundation and kernel checker, we also have  
 906 to trust that code extraction respects that foundation, that the implementation of the fold  
 907 operation for hashmaps satisfies its description, and that, after parsing, the computation  
 908 uses the correct input formula. All of our specifications have been completely proven and  
 909 verified relative to these assumptions.

## 910 **D Detailed Experimental Results**

911 Our experiments are designed to evaluate the following:

- 912 ■ whether our toolchain can handle challenging benchmark problems,
- 913 ■ the effectiveness of some design choices and optimizations,
- 914 ■ the comparative performance of the prototype and verified checkers, and
- 915 ■ how the performance of our toolchain compares to that of the MICE model counter  
 916 verifier.

### 917 **D.1 Experimental Setup**

918 As a test machine, we used a 2021 MacBook Pro laptop having a 3.2 Ghz Apple M1 processor  
 919 and with 64 GB of physical memory. We also used an external 3 TB solid-state disk drive  
 920 with an advertised read/write speed of one gigabit per second. Despite being a laptop, this  
 921 configuration is somewhat more powerful than the nodes typically found in server clusters.  
 922 The fast access to storage is especially important due to the very large files generated and  
 923 processed by the tools. Our tools generated individual files with over 160 gigabytes of data.

924 For benchmarks, we downloaded 100 files each from the 2022 standard and weighted  
 925 model counting competitions:

926 `https://mccompetition.org/2022/mc\_description.html`

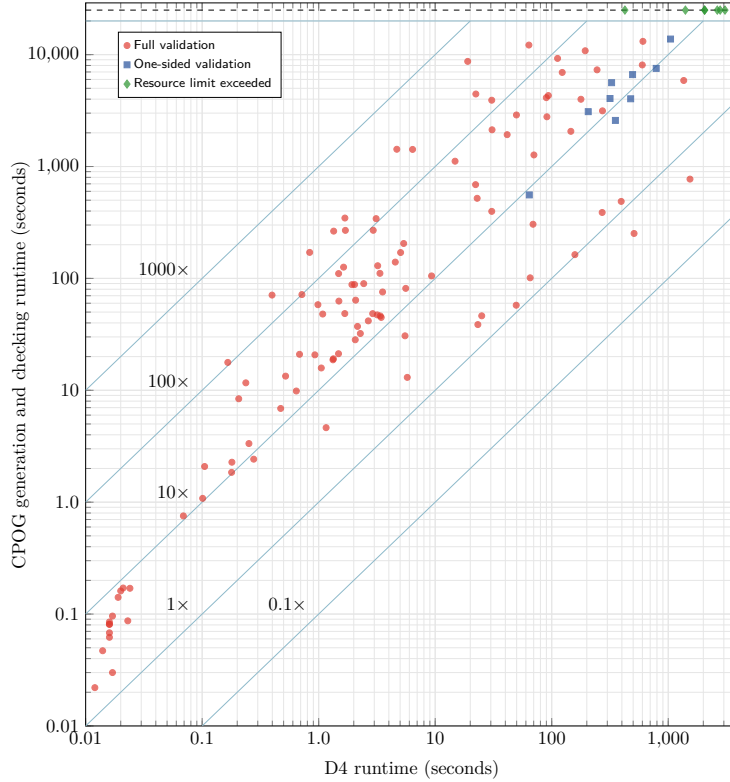
927 We found that 20 of these were duplicates across the two tracks, yielding a total of 180  
 928 unique benchmark problems, ranging in size from 250 to 2,753,207 clauses.

929 Our default configuration for the proof generator used the structural method for the  
 930 forward implication proof, with the two optimizations of literal grouping and lemmas. Running  
 931 with a time limit of 4000 seconds<sup>2</sup>, D4 was able to generate d-DNNF representations for 124  
 932 of these, and it timed out on the rest. Our proof generator was able to convert all of the

---

<sup>2</sup> All measured times in this document are actual elapsed times, not CPU times.





■ **Figure 4** Combined runtime for CPOG proof generation and checking as function of D4 runtime. Timeouts are shown as points on the dashed line.

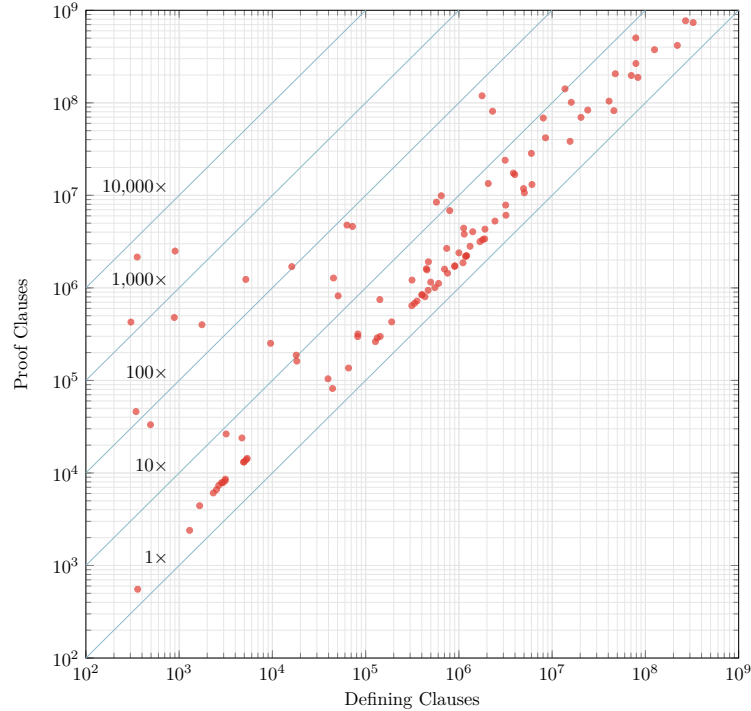
933 generated d-DNNF graphs into POGs and determine how many defining clauses they would  
 934 generate. A POG operation with  $k$  arguments requires  $k + 1$  defining clauses, and so the  
 935 total number of defining clauses in POG  $P$  equals the number of nonterminal nodes plus  
 936 the number of edges in the graph, corresponding to our definition of  $|P|$  in Section 5. The  
 937 number of defining clauses ranged from 304 to 2,761,457,765 with a median of 848,784.

## 938 D.2 Performance of the Proof Generator and Checker

939 For each of the 124 d-DNNF graphs generated by D4, we ran our proof generator and  
 940 prototype checker, limiting proof generation to 10,000 seconds. We ran the programs to  
 941 generate and check one-sided proofs for the graphs, again with a time limit of 10,000 seconds  
 942 for proof generation.

943 Figure 4 summarizes our results in terms of the time required by D4 (X axis) versus the  
 944 sum of the times for the proof generator and checker (Y axis). We were able to complete a  
 945 full validation for 108 of the 124 benchmarks, with times ranging from fractions of a second  
 946 to 13,144 seconds, with a median of 71.6 seconds.

947 Relative to the runtime for D4, two problems ran faster with the generator and checker.  
 948 These were for problems having small numbers of models (relative to the number of variables),  
 949 and so most of the time spent by both D4 and our proof generator was in running a CDCL  
 950 solver to search the very sparse solution space. CADICAL generally outperforms the miniSAT-  
 951 based solver used by D4. At the other extreme, one problem that required only 19 seconds  
 952 for D4 required 8657 seconds to generate a proof and 45 seconds to check it. This benchmark



■ **Figure 5** Total number of clauses in CPOG file as function of number of defining clauses

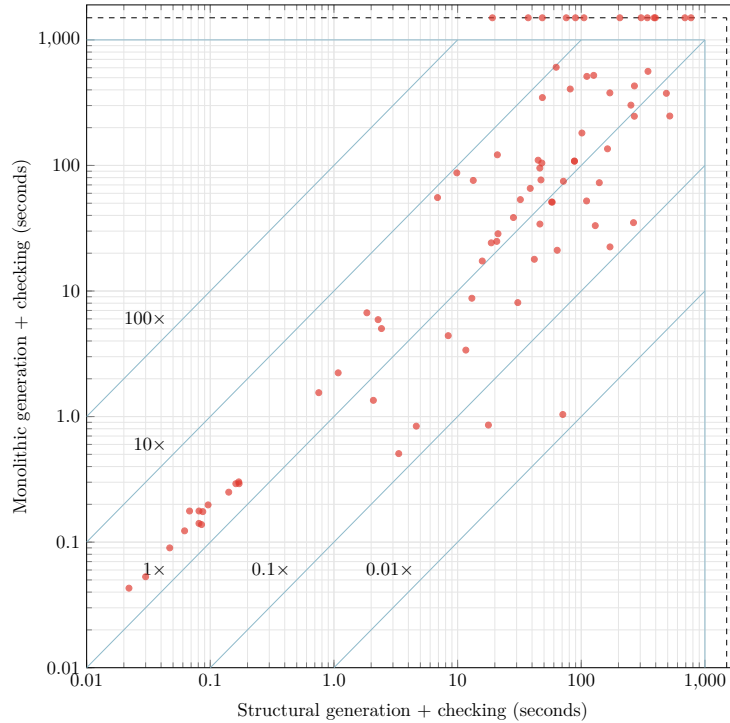
953 appears to stem from weak performance by our implementation of BCP. Overall the ratios  
 954 between the combined generation plus checking times versus the time for D4 had a harmonic  
 955 mean of 5.5.

956 Of the 16 benchmarks that could not be fully validated, one had so many defining clauses  
 957 that it overflowed the 32-bit signed integers our programs use for clause identifiers. Another  
 958 exited due to the virtual memory limit imposed by the operating system, and the other 14  
 959 timed out. Of the 15 that did not overflow the clause counter, we were able to complete a  
 960 one-sided verification for 9, but the other 6 timed out. Overall, we were able to provide some  
 961 level of verification for 94% of the benchmark problems.

962 Figure 5 compares the number of defining clauses (X axis) with the total number of  
 963 clauses (Y axis) for the 108 problems that were fully verified. These totals include the  
 964 defining clauses for the POG, the additional defining clauses introduced to support literal  
 965 grouping and lemmas, and the clauses added by RUP steps in the forward implication proof.  
 966 These totals ranged from 554 to 770,382,773, with a median of 1,719,245. The ratios of the  
 967 total clauses versus defining clauses ranged from 1.54 to 6073, with a harmonic mean of  
 968 3.13. The high numbers were for problems having very few models, and so the proof clauses  
 969 must encode large proofs of unsatisfiability.

### 970 D.3 Monolithic Forward Implication Proofs

971 To assess the relative merits of the two approaches to forward implication proof generation,  
 972 we ran the proof generator in monolithic mode on the 92 problems for which the combination  
 973 of structural proof generation and checking required at most 3000 seconds and then kept  
 974 only those measurements for which at least one of the two approaches had combined times  
 975 (generation plus checking) of at most 1000 seconds. There were a total of 83 such problems.



■ **Figure 6** Monolithic versus Structural Validation: Time. Timeouts are shown as points on the dashed lines.

976 Figures 6 and 7 show comparisons between the two approaches in terms of time and total  
 977 clauses.

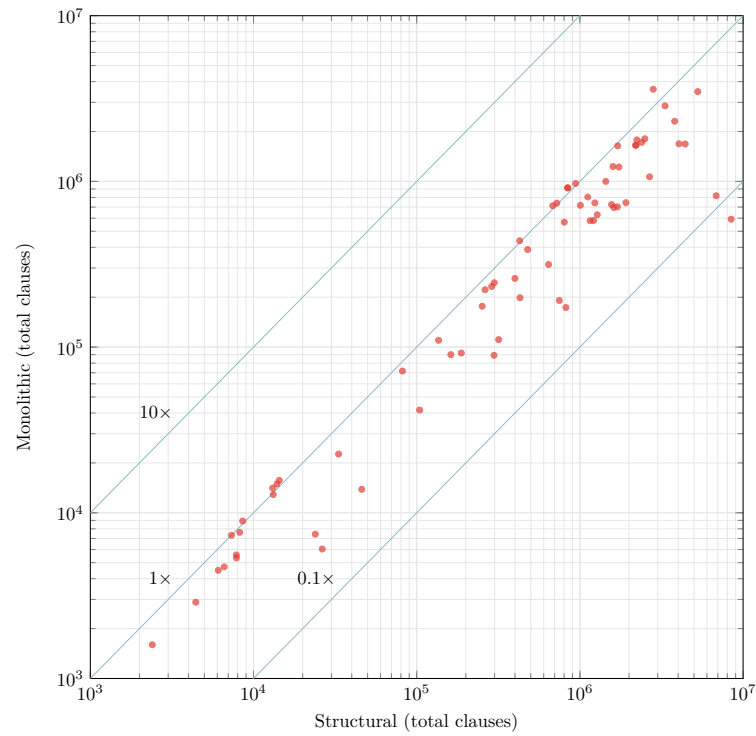
978 Examining Figure 6, we can see that the monolithic approach ran faster than the structural  
 979 approach for 23 of the problems, including a case where the generator and checker required  
 980 only one second, versus 71 seconds for the structural approach. On the other hand, the  
 981 monolithic approach timed out for 13 of the cases and was slower for the other 37. In general,  
 982 we can see the monolithic approach faring worse on larger problems.

983 Figure 7 shows the total number of clauses for the two approaches for the 70 cases where  
 984 both approaches completed within 1000 seconds. As can be seen, the monolithic approach  
 985 consistently produces shorter proofs, perhaps due to the benefits of the proof trimming by  
 986 DRAT-TRIM.

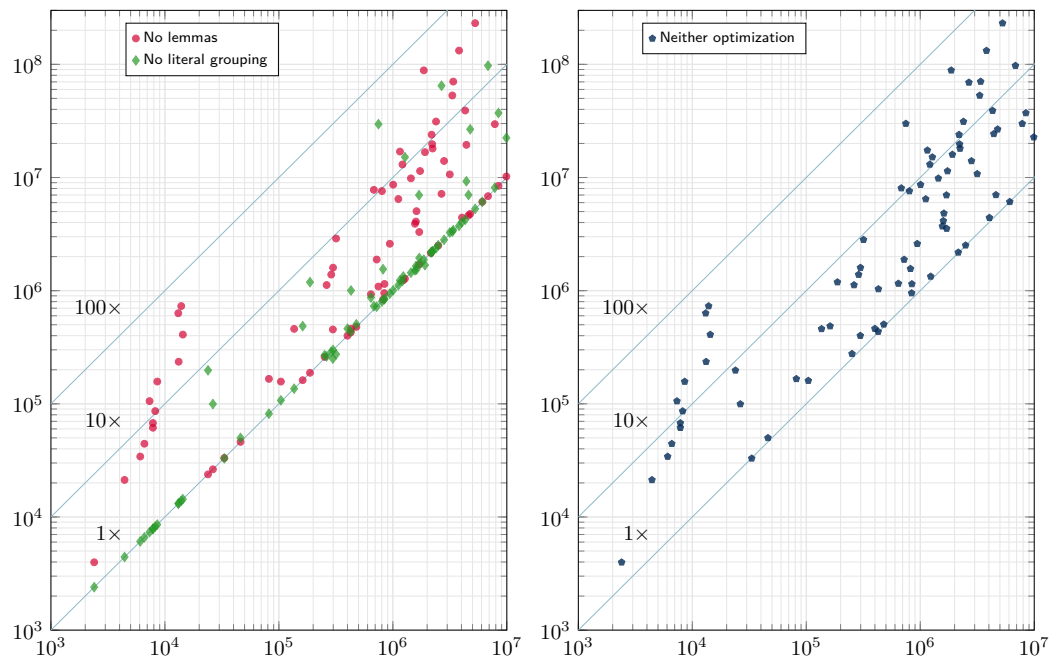
987 These experiments suggest that a hybrid of the two approaches might yield the best  
 988 results in terms of consistency, runtime, and proof size. It would begin the recursive descent  
 989 of `validate` according to a structural approach, but monolithically generate the proof once it  
 990 encounters a sufficiently small subgraph.

## 991 D.4 Impact of Optimizations

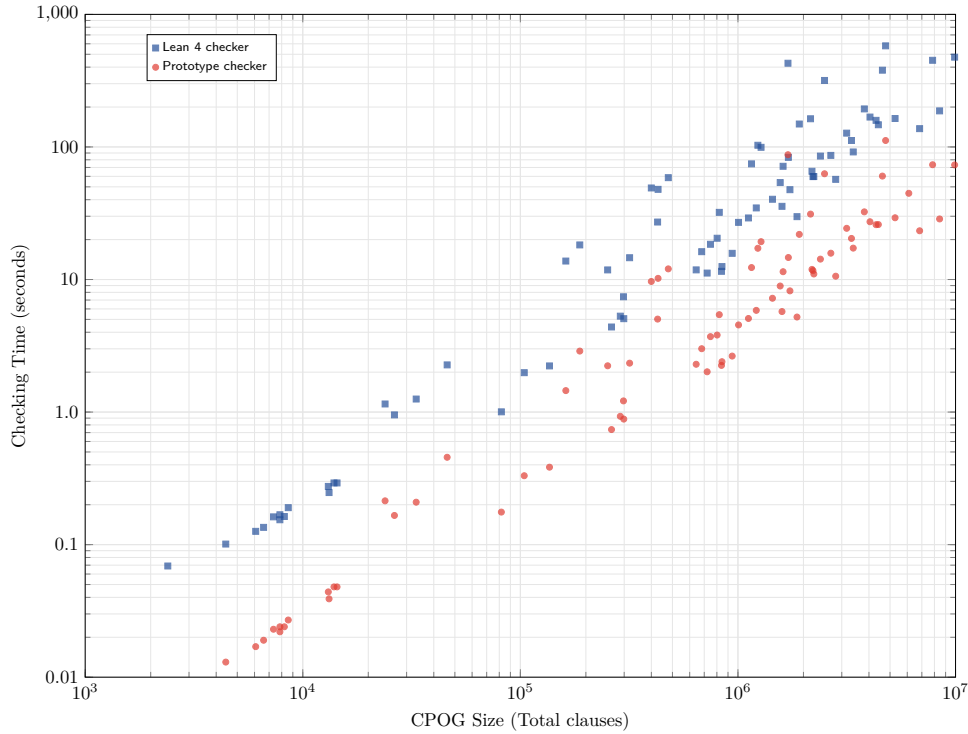
992 To assess the impact of the two optimizations: literal grouping and lemmas, we ran the  
 993 proof generator on the 80 problems for which the combination of the two optimizations  
 994 yielded proofs totaling at most 10 million clauses. Compared to the problems used to assess  
 995 monolithic proof generation (Section D.3), this set included one that had a short proof but  
 996 long runtime and excluded four that had short runtimes but long proofs. We then ran the  
 997 proof generator on all three combinations of the optimizations being partially or totally



■ **Figure 7** Monolithic versus Structural Validation: Clauses



■ **Figure 8** Proof clauses when one (left) or both (right) optimizations are disabled, versus when both optimizations are enabled



■ **Figure 9** Times for Prototype and Verified Checker, as Functions of Total Clauses

998 disabled. Figure 8 shows how disabling these optimizations affected the total number of  
 999 clauses in the proofs, where the X axis indicates the total number with both optimizations  
 1000 enabled, while the Y axis indicates the number with one (left) or both (right) optimiza-  
 1001 tions disabled.

1002 As the figures illustrated, each optimization on its own can shorten the proofs, sometimes  
 1003 substantially, and the two combine to have even greater effect. Of the 80 benchmark problems  
 1004 evaluated, and comparing with both optimizations disabled, 47 had the number of proof  
 1005 clauses reduced by a factor of 2 by using lemmas, 14 by using literal merging, and 60 by  
 1006 enabling both optimizations. In the extreme cases, a lack of lemmas caused one proof to  
 1007 grow by a factor of 52.5, while a lack of literal grouping causes another proof to grow by a  
 1008 factor of 39.6.

1009 The time performance of these optimizations is less dramatic, but still significant. Lemmas  
 1010 sometimes made proof generation slower, but in one case ran  $38\times$  faster. Similarly, literal  
 1011 grouping could take longer, but in one case ran  $3.9\times$  faster.

1012 Overall, it is clear that these two optimizations are worthwhile, and sometimes critical,  
 1013 to success for some benchmarks.

## 1014 D.5 Evaluation of the Verified Checker

1015 We ran the Lean 4 checker on the 80 benchmark problems for which the generated CPOG  
 1016 file had at most 10 million total clauses. These are the same benchmarks used in evaluating  
 1017 the optimizations (Section D.4). All of the checks completed and confirmed the outcome of  
 1018 the prototype checker.

1019 Figure 9 shows the performance of both the Lean 4 checker and the prototype checker on

1020 the Y axis, as functions of the total number of clauses on the X axis. As would be expected,  
1021 the Lean 4 checker consistently runs slower than the prototype checker. The ratio of the  
1022 runtime for Lean 4 versus the runtime for the prototype checker had a harmonic mean of  
1023 5.94.

1024 Importantly, however, it can be seen that both checkers show the same overall performance  
1025 trends. The ratio of the two runtimes was less than 8.0 for all but three small benchmarks.  
1026 This seems like a reasonable price to pay for a rigorous guarantee of correctness, and we are  
1027 confident that we can reduce this gap with more effort in optimizing the verified checker.

## 1028 **D.6 Comparison with the MICE Framework**

1029 To compare the performance of our toolchain versus that of MICE [8], we consider how the  
1030 combination of proof generation, proof checking, and ring evaluation in the CPOG toolchain  
1031 compares to the time to run their proof generator `NNF2TRACE` and checker `SHARPTRACE`. We  
1032 add the time for ring evaluation on our side, since their checker provides a count, although  
1033 only an unweighted one. We used the prototype versions of the checker and ring evaluator.  
1034 We ran their tools on the 92 benchmark problems for which our toolchain requires at most  
1035 3000 seconds and retained the 84 data points for which at least one of the toolchains completes  
1036 in under 1000 seconds.

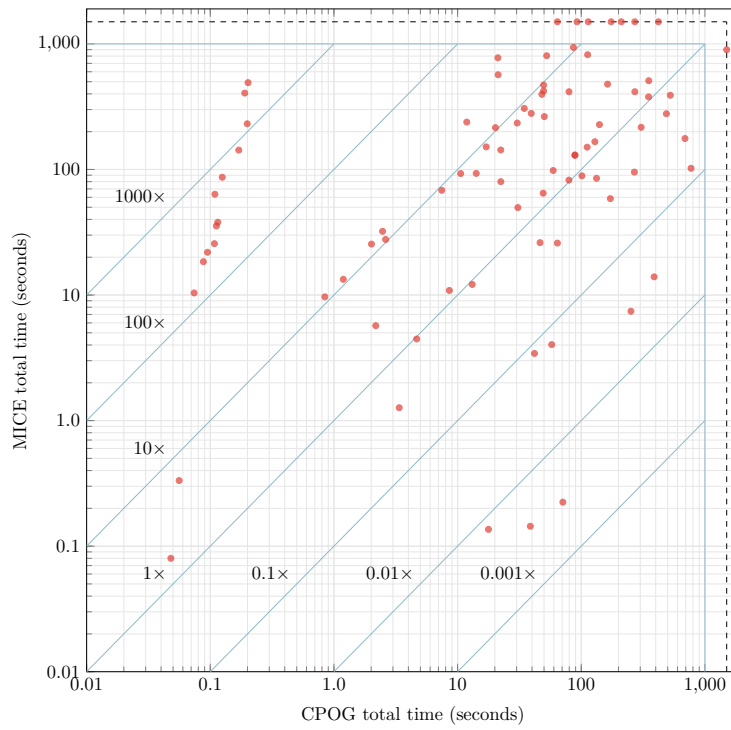
1037 We can summarize the results as:

- 1038 ■ 84 total data points
- 1039 ■ 76 were completed by both programs in under 1000 seconds
- 1040 ■ Of those 21 were faster with MICE, 55 with CPOG
- 1041 ■ 7 completed with CPOG under 1000 seconds but exceeded that threshold for MICE
- 1042 ■ 1 completed with MICE under 1000 seconds but exceeded that threshold for CPOG

1043 The detailed results are shown in Figure 10. As can be seen, the comparative runtimes  
1044 are highly variable, reflecting the fact that the two toolchains differ fundamentally in their  
1045 objectives and their approaches.

1046 One shortcoming of `NNF2TRACE` can be seen in the near-vertical series of points in the  
1047 upper-lefthand corner. These correlate with a similar set of points along the left in Figure 8.  
1048 Like the naive implementation of `validate`, their program recursively traverses the graph  
1049 representation of a formula, effectively expanding it into a tree. They lack an optimization  
1050 analogous to lemmas.

1051 We have not tried the MICE toolchain on larger benchmarks, but we suspect it would  
1052 encounter significant scaling limitations.



■ **Figure 10** Running Time for MICE versus CPOG proof chains. Timeouts are shown as points on the dashed lines.

