Boosting Constrained Horn Solving by Unsat Core Learning

Parosh Aziz Abdulla^{*1}, Chencheng Liang^{*1}, and Philipp Rümmer^{*1,2}

¹ Uppsala University, Uppsala, Sweden chencheng.liang@it.uu.se parosh.abdulla@it.uu.se
² University of Regensburg, Regensburg, Germany philipp.ruemmer@ur.de

Abstract. The Relational Hyper-Graph Neural Network (R-HyGNN) was introduced in [1] to learn domain-specific knowledge from program verification problems encoded in Constrained Horn Clauses (CHCs). It exhibits high accuracy in predicting the occurrence of CHCs in counterexamples. In this research, we present an R-HyGNN-based framework called MUSHyperNet. The goal is to predict the Minimal Unsatisfiable Subsets (MUSes) (i.e., unsat core) of a set of CHCs to guide an abstract symbolic model checking algorithm. In MUSHyperNet, we can predict the MUSes once and use them in different instances of the abstract symbolic model checking algorithm. We demonstrate the efficacy of MUSHyperNet using two instances of the abstract symbolic modelchecking algorithm: Counter-Example Guided Abstraction Refinement (CEGAR) and symbolic model-checking-based (SymEx) algorithms. Our framework enhances performance on a uniform selection of benchmarks across all categories from CHC-COMP, solving more problems (6.1% increase for SymEx, 4.1% for CEGAR) and reducing average solving time (13.3% for SymEx, 7.1% for CEGAR).

Keywords: Automatic program verification · Constrained Horn clauses · Graph Neural Networks.

1 Introduction

Constrained Horn Clauses (CHCs) [2] are logical formulas that can describe program behaviors and specifications. Encoding program verification problems in CHCs and solving them (checking CHCs' satisfiability) has been an active research area for a number of years [3,4]. If the encoded CHCs are satisfiable, the corresponding program verification problem is safe; if not, it is unsafe.

Solving a set of CHCs means that we either find an interpretation for the predicate (relation) symbols and variables that satisfies all the clauses, or prove

^{*} Author names in alphabetical order. The theory presented in the paper was developed by Abdulla, Liang, and Rümmer, the implementation is by Liang and Rümmer, evaluation was done by Liang.

that no such interpretation exists. Various techniques, such as Counterexample Guided Abstraction Refinement (CEGAR) [5] and IC3 [6], have been utilized for this purpose. However, due to the undecidability of solving CHCs, we need carefully designed or tuned heuristics for specific instances.

In this paper, we consider Minimal Unsatisfiable Subsets (MUSes) [7] of sets of CHCs to support the solving process. Given an unsatisfiable set of CHCs, each MUS is a subset that is again unsatisfiable, but removing any CHC from an MUS makes it satisfiable. Understanding MUSes of a set of CHCs can guide solvers to focus on error-prone clauses [8, Section 3.1]: for an unsatisfiable set of CHCs, evaluating the satisfiability of MUSes first can quickly identify unsatisfiable CHCs, eliminating the need for a comprehensive check. In a satisfiable set of CHCs, examining MUSes first can provide a better starting point for refining potentially problematic constraints. This guidance can help the solver converge towards a solution more efficiently. Manually designed heuristics to find MUSes involves summarizing and generalizing the features of a set of CHCs from examples, which can be replaced by data-driven methods.

Various studies that apply deep learning to formal methods for verification have been published in the recent past, e.g., [9,10,11,12]. Primarily, deep learning serves as a feature extractor, which can automatically summarize and generalize program features from examples, alleviating the need for manual crafting and tuning of various heuristics. In addition, the idea of representing logic formulas as graphs and using Graph Neural Networks (GNNs) [13] to guide solvers has been employed in many successful studies [14,15,16,17].

However, we are not aware of any study that applies GNNs to guide CHCbased symbolic model checking techniques. The main contribution of this paper is to train a GNN model to predict MUSes of a set of CHCs. We train a GNN using the CHC-R-HyGNN framework [1] to predict values between 0 and 1, representing the probabilities of CHCs being elements of MUSes. For example, we assume that a set $C = \{c_1, c_2, c_3\}$ of CHCs has one MUS $\{c_1, c_2\}$. The predicted probabilities for c_1, c_2 , and c_3 being in the MUS could be 0.9, 0.8, and 0.1, respectively. We propose several strategies that use the predicted probability to guide an abstract symbolic model checking algorithm. The strategies can be instantiated for different algorithms on CHCs, such as CEGAR- and symbolic execution-based satisfiability checkers.

Figure 1 depicts an overview of our framework MUSHyperNet. Firstly, we encode the program verification problem into a set of CHCs. Secondly, we use the graph encoder in the CHC-R-HyGNN framework [1] to convert the set of CHCs into a graph format. Then, we train a GNN model named Relational Hypergraph Neural Network (R-HyGNN) to predict the probability of each CHC occurring in MUSes. Finally, we employ the predicted probabilities to guide the abstract symbolic model checking algorithm by determining the sequence for processing each CHC in the set of CHCs.

We utilize the same benchmarks as in [1], comprising 17130 Linear Integer Arithmetic (LIA) problems from the CHC-COMP 2021 repository [18] for training and evaluating our approach. Further details about the benchmark can



Fig. 1: The CHC-R-HyGNN framework [1] (represented by the round box) comprises a CHC graph encoder and a GNN known as the Relational Hypergraph Neural Network (R-HyGNN), which is capable of handling hypergraphs. In our previous work, we introduced the CHC-R-HyGNN framework that employs various proxy tasks to generalize valuable information for constructing heuristics for CHC-encoded program verification problems.

be found in [19, Table 1]. The problems for evaluation are uniformly selected from this benchmark and can be found in a public repository [20] The experimental results show an improvement of up to 4.1% and 6.1% in the number of solved problems for the CEGAR and SymEx algorithms, respectively. Additionally, the average solving time demonstrates enhancements of 7.1% and 13.3% for the CEGAR and SymEx algorithms, respectively. In other words, MUSHyperNet can increase the number of solved problems and decrease the solving time for problems similar to those in the CHC-COMP benchmark. To the best of our knowledge, this is the first time unsat core learning has been used successfully in the context of CHCs.

In summary, our contributions are as follows:

- We develop a GNN-based framework named MUSHyperNet, which trains a GNN to predict the MUSes of CHCs and utilizes the predicted probabilities to guide an abstract symbolic model checking algorithm.
- We explore GNN models trained on different datasets and methods for applying predicted MUSes to guide two instances of the model checking algorithm.
- We evaluate MUSHyperNet on 383 linear and 488 non-linear LIA problems, uniformly sampled from the CHC-COMP benchmark [19]. The improvements in the number of solved problems and average solving time are up to 6.1% and 13.3% for the SymEx, and 4.1% and 7.1% for the CEGAR.

2 Preliminaries

We first introduce required notation for multi-sorted first-order logic, and define Constrained Horn Clauses (CHCs) and the encoding of a program verification problem as CHCs. Finally, we explain basic concepts of Graph Neural Network (GNN) and introduce Relational Hyper Graph Neural Network (R-HyGNN).

2.1 Notations

We assume familiarity with standard multi-sorted first-order logic (e.g., see [21]). A first-order language \mathcal{L} is defined by a signature $\mathcal{L} = (\mathcal{S}, \mathcal{R}, \mathcal{F}, \mathcal{X})$, where \mathcal{S} is

a non-empty set of sorts; \mathcal{R} is a set of fixed-arity predicate (relation) symbols, each of which is associated with a list of argument sorts; \mathcal{F} is a set of fixedarity function symbols, each of which is associated with a list of argument sorts and a result sort; and $\mathcal{X} = \bigcup_{s \in \mathcal{S}} \mathcal{X}_s$ is a set of sorted variables, where \mathcal{X}_s are the variables of sort s. A term t is a variable from \mathcal{X} , or an n-ary function symbol $f \in \mathcal{F}$ applied to terms t_1, \ldots, t_n of the right sorts. An *atomic formula* (*atom*) of \mathcal{L} is of the form $p(t_1, \ldots, t_n)$, where $p \in \mathcal{R}$ is an n-ary predicate symbol and t_1, \ldots, t_n are terms of the right sorts. A *formula* is a Boolean literal *true*, *false*, an atom, or obtained by applying logical connectives \neg , \land , \lor , \rightarrow and quantifiers \forall , \exists to formulas. We write implications both left-to-right ($\varphi \rightarrow \psi$) and right-to-left ($\psi \leftarrow \varphi$). A formula is *closed* if all variables occurring in the formula are bound by quantifiers.

A multi-sorted structure $\mathcal{M} = (\mathcal{U}, I)$ for \mathcal{L} consists of a set $\mathcal{U} = \bigcup_{s \in S} \mathcal{U}_s$, being the union of non-empty domains \mathcal{U}_s of each sort $s \in S$, and an interpretation I such that $I(s) = \mathcal{U}_s$ for every sort $s \in S$; for each n-ary predicate symbol $p \in \mathcal{R}$ with argument sorts s_1, \ldots, s_n , $I(p) \subseteq U_{s_1} \times \cdots \times U_{s_n}$; and for each n-ary function symbol $f \in \mathcal{F}$ with argument sorts s_1, \ldots, s_n and result sort s, $I(f) \in U_{s_1} \times \cdots \times U_{s_n} \to U_s$. A variable assignment β for the structure $\mathcal{M} = (\mathcal{U}, I)$ is a function $\mathcal{X} \to \mathcal{U}$ that maps each variable $x \in \mathcal{X}_s$ to an element of the corresponding domain \mathcal{U}_s . Given \mathcal{L} , a structure $\mathcal{M} = (\mathcal{U}, I)$, and a variable assignment β , the evaluation of a term or formula is performed by the function $val_{\mathcal{M},\beta}$, defined by $val_{\mathcal{M},\beta}(x) = \beta(x)$ for a variable $x \in \mathcal{X}$; $val_{\mathcal{M},\beta}(f(t_1,\ldots,t_n)) = I(f)[val_{\mathcal{M},\beta}(t_1),\ldots,val_{\mathcal{M},\beta}(t_n)]$ for a function $f \in \mathcal{F}$; and $val_{\mathcal{M},\beta}(p(t_1,\ldots,t_n)) = true$ iff $(val_{\mathcal{M},\beta}(t_1),\ldots,val_{\mathcal{M},\beta}(t_n)) \in I(p)$. The evaluation of compound formulas is defined as is common. When \mathcal{M} is clear from the context, we also write val_{β} instead of $val_{\mathcal{M},\beta}$.

We say that a formula φ is *satisfied* in \mathcal{M}, β if $val_{\mathcal{M},\beta}(\varphi) = true$, and that it is *satisfiable* (SAT) if it is satisfied by some \mathcal{M}, β . We say a set Γ of formulae *entails* a formula φ , denoted $\Gamma \models \varphi$, if φ is satisfied whenever all formulas in Γ are satisfied.

Example 1. We assume a language \mathcal{L} consisting of a sort s, two constants a and b, a unary function symbol f with argument and result of sort s; a binary predicate symbol p with two arguments of sort s. A structure $\mathcal{M} = (U_s, I)$ can be defined by $U_s = \mathbb{Z}$, I(a) = 1, I(b) = 2, I(f)[x] = x + 2, and $I(p)[x, y] = x \leq y$. The formula $\varphi = p(a, b) \rightarrow p(x, f(a))$ is satisfied in \mathcal{M} by a variable assignment $\beta(x) = 1$ since $val_{\beta}(\varphi) = I(p)[val_{\beta}(a), val_{\beta}(b)] \rightarrow I(p)[val_{\beta}(x), I(f)[val_{\beta}(x))] = \neg(1 \leq 2) \lor 1 \leq (1+2)$ is true. The formula φ is satisfiable since $val_{\beta}(\varphi) = true$ for \mathcal{M} and a variable assignment $\beta(x) = 1$.

2.2 Constraint Horn Clauses

To introduce the notion of constrained Horn clauses, we assume a fixed base signature $\Sigma = (S, \mathcal{R}, \mathcal{F}, \mathcal{X})$, as well as a unique structure \mathcal{M} over this signature, forming the background theory. In this paper, we mainly consider the background theory of Linear Integer Arithmetic (LIA), following the SMT-LIB standard [22]. We further assume a set \mathcal{R}_C of additional relation symbols that is disjoint from \mathcal{R} , which will be used to formulate the head and body of clauses.

Definition 1 (Constrained Horn clause). Given signature Σ and the set \mathcal{R}_C , a Constrained Horn Clause (CHC) is a closed formula in the form

$$\forall \bar{x}. \ H \leftarrow p_1(\bar{t}_1) \land \dots \land p_n(\bar{t}_n) \land \varphi, \tag{1}$$

where \bar{x} is a vector of variables; H is either false or an atom $p(t_1, \ldots, t_n)$ with $p \in \mathcal{R}_C$; the relation symbols p_1, \ldots, p_n are elements of \mathcal{R}_C ; and $\bar{t}_1, \ldots, \bar{t}_n$ and φ are vectors of terms and a formula over Σ , respectively. We call H the head and $p_1(\bar{t}_1) \wedge \cdots \wedge p_n(\bar{t}_n) \wedge \varphi$ the body of the clause, respectively. We call the formula φ in the body the constraint of the clause.

For convenience, in many places we leave out the quantifiers $\forall \bar{x} \text{ when writing clauses.}$ A CHC without atoms in its body (the case n = 0) is called a *fact*. If the body of a CHC contains zero or one atom, the CHC is called *linear*. Otherwise, it is called *non-linear*.

Solving CHCs boils down to searching interpretations of the relation symbols \mathcal{R}_C that satisfy the CHCs, assuming that all background symbols from Σ are interpreted by the fixed structure \mathcal{M} :

Definition 2 (Satisfiability of a CHC). A CHC h is satisfiable if there is a structure $\mathcal{M}_C = (\mathcal{U}, I_C)$ for the extended signature $\Sigma_C = (\mathcal{S}, \mathcal{R} \cup \mathcal{R}_C, \mathcal{F}, \mathcal{X})$ such that (i) I_C coincides with I on Σ , and (ii) \mathcal{M}_C satisfies h. A set C of CHCs is satisfiable if there is an extended structure \mathcal{M}_C simultaneously satisfying all clauses in C.

Encoding Program Verification Problems using CHCs. A program verification problem involves checking whether a program adheres to its specified behavior. One approach to verification is to transform the problem into determining the satisfiability of a set of CHCs. This can be done, for instance, by encoding the partial correctness of a procedural imperative program into a negated Existential positive Least Fixed-point Logic (E+LFP) formula [4] using the weakest precondition calculus. Generally, encodings are designed such that the set of CHCs is satisfiable if and only if a program is safe. Various encoding schemes for different programming languages have been introduced in the literature, e.g., [4].

2.3 Graph Neural Networks

A Graph Neural Network (GNN) [13] is a type of neural network that consists of Multi-Layer Perceptrons (MLPs) [23]. GNNs operate on graph-structured data with nodes and edges, making them suitable for logic formulas which can naturally be represented as graphs. A GNN can take a set of typed nodes and edges as input and output a set of feature representations (vectors of real numbers) associated with the properties of the nodes. We refer to them as *node representations* in the rest of the sections.

The Message-Passing based GNN (MP-GNN) [24] is a type of GNN model. It utilizes an iterative message-passing algorithm in which each node in the graph aggregates messages from its neighboring nodes to update its own node representation. This mechanism assists in identifying the inner connections within substructures, such as terms and atoms, in graph represented logic formulae.

Formally, let G = (V, E) be a graph, where V is the set of nodes and E is the set of edges. Let x_v be the initial node representation (a vector of random real numbers) for node v in the graph, and let N_v be the set of neighbors of node v. An MP-GNN consists of a series of T message-passing steps. At each step t, every node v in the graph updates its node representation as follows:

$$h_v^t = \phi_t(\rho_t(\{h_u^{t-1} \mid u \in N_v\}), h_v^{t-1}),$$
(2)

where $h_v^t \in \mathbb{R}^n$ is the updated node representation for node v after t iterations. The initial node representation, h_v^0 , is usually derived from the node type and given by x_v . The node representation of u in the previous iteration t-1 is h_u^{t-1} , and node u is a neighbor of node v. $\rho_t : (\mathbb{R}^n)^{|N_v|} \to \mathbb{R}^n$ is a aggregation function with trainable parameters (e.g., a MLP followed by sum, min, or max) that aggregates the node representations of v's neighboring nodes at the t-th iteration. $\phi_t : \mathbb{R}^n \to \mathbb{R}^n$ is a function with trainable parameters (e.g., a MLP) that takes the aggregated node representation from ρ_t and the node representation of vin previous iteration as input, and outputs the updated node representation of v at the t-th iteration. MP-GNN assumes a node can capture local structural information from t-hop's neighbors by updating the node representation using aggregated representations of the neighbor nodes.

The final output of the MP-GNN could be the set of updated node representations for all nodes in the graph after T iterations. These node representations can be used for a variety of downstream tasks, such as node classification or graph classification.

Relational Hyper-Graph Neural Network (R-HyGNN) [1] is an extension of one MP-GNN called Relational Graph Convolutional Networks (R-GCN) [25], and it is specifically designed to handle labeled hypergraphs.

A labeled (typed) hypergraph is a hypergraph where each vertex (node) and hyperedge is assigned a type from a predefined set of types. Formally, a labeled hypergraph LHG is defined as a tuple $LHG = (V, E, \lambda_V, L_V, L_E)$, where V is a set of elements called vertices (or nodes), L_E is a set of pair consisting of a label (type) r and the number of nodes k under the label $r, E \subseteq V^* \times L_E$ is a set of hyperedges in which each hyperedge consists of a non-empty subsets of V and a pair $(r, k) \in L_E$. Here, $\lambda_V : V \to L_V$ is a labeling function that assigns a type from the set L_V to each vertex in V. The L_V is a set of possible types (labels) for the vertices V.

The node representation updating rule of R-HyGNN for one node v at timestep t is

$$h_{v}^{t} = \operatorname{ReLU}(\sum_{\substack{(w_{1},\dots,w_{k},(r,k))\\\in E}}\sum_{i\in\{1,\dots,k\},\\w_{i}=v}W_{r,i}^{t} \cdot \|(h_{w_{1}}^{t-1},\dots,h_{w_{i-1}}^{t-1},h_{w_{i+1}}^{t-1},\dots,h_{w_{k}}^{t-1})),$$
(3)

 $\mathbf{6}$

where the pair $(r, k) \in L_E$ is the edge type (relation) and the number of node for a edge $(w_1, \ldots, w_k, (r, k)) \in E$, $W_{r,i}^t$ is a matrix of learnable parameters in time step t for node $v = w_i$ in the edge with type r. There are $|L_E| \times \sum_{(r,k) \in L_E} (k) \times t$ matrices of learnable parameters in total. Here, $\|(h_{w_1}^{t-1}, \ldots, h_{w_{i-1}}^{t-1}, h_{w_{i+1}}^{t-1}, \ldots, h_{w_k}^{t-1})\|$ means concatenate v' neighbour node representations in time step t-1. The initial node representation h_v^0 is derived from the node types L_V .

Intuitively, to update the representation of a node, R-HyGNN first concatenates the neighbor representations of node v for each edge from the previous time step t-1. It then multiplies the concatenated neighbor representations by the corresponding matrix of trained parameters (i.e., $W_{r,i}^t$) to derive a local representation of v. Next, it aggregates the local node representations (e.g., by addition). In other words, $\sum_{\substack{(w_1,\ldots,w_k,(r,k))\\\in E}} \sum_{i\in\{1,\ldots,k\}}, W_{r,i}^t \cdot \|(h_{w_1}^{t-1},\ldots,h_{w_{i-1}}^{t-1},h_{w_{i+1}}^{t-1},\ldots,h_{w_k}^{t-1})$ can be abstract to $\rho_t(\|(h_{w_1}^{t-1},\ldots,h_{w_{i-1}}^{t-1},h_{w_{i+1}}^{t-1},\ldots,h_{w_k}^{t-1}))$ where ρ_t is an aggregation function with trainable parameters in $W_{r,i}^t$. Finally, it applies a ReLU function [26] as the update function ϕ_t . This function takes the aggregated local node representations as input and produces the final node representation h_v^t . This updating process for a single node recurs t times.

3 Abstract Symbolic Model Checking for CHCs

The goal of our work is to utilize GNNs to obtain improved state space exploration methods for CHCs. To this end, in this section we introduce an abstract formulation of CHC state space exploration, covering both the classical CEGAR approach and exploration in the style of symbolic execution.

3.1 Satisfiability Checking for CHCs

Our Algorithm 1 checks the satisfiability of a given a set C of CHCs by constructing an abstract reachability hyper-graph (ARG). An ARG is an overapproximation of the facts $p(\bar{a})$ that are logically entailed by C; by demonstrating that the atom *false* does not follow from C, it can be shown that C is satisfiable. Since each node of an ARG can represent a whole set of facts, a finite ARG can be a representation even of infinite models of a set C.

We first give an abstract definition of ARGs that does not mandate any particular symbolic representation of sets of $p(\bar{a})$. We later introduce two instances of this abstract framework.

Like in Section 2.2, we denote the set of relation symbols used in CHCs by \mathcal{R}_C . For a k-ary relation symbol $p \in \mathcal{R}_C$ with argument sorts s_1, \ldots, s_k , we write $\mathcal{R}_p = \mathcal{P}(U_{s_1} \times \cdots \times U_{s_k})$ for the set of possible relations represented by p.

Definition 3. An abstract reachability graph for a set C of CHCs is a hypergraph (V, E), where

- the set V of nodes is a set of pairs (p, R), with p being a relation symbol and $R \in \mathcal{R}_p$;

 $-E \subseteq V^* \times \mathcal{C} \times V$ is a set of hyper-edges labelled with clauses. For every edge $(v_1, \ldots, v_n, h, v_0) \in E$, it is the case that:

- the head of a clause h is not false, i.e., h is of the form $\forall \bar{x}. p_0(\bar{t}_0) \leftarrow p_1(\bar{t}_1) \wedge \cdots \wedge p_n(\bar{t}_n) \wedge \varphi;$
- nodes v_0, \ldots, v_n correspond to the head and the body of h, i.e., for $i \in \{0, \ldots, n\}$ it is the case that $v_i = (p_i, R_i)$ for some $R_i \in \mathcal{R}_{p_i}$;
- R_0 over-approximates the facts implied by the clause h:

$$R_0 \supseteq \left\{ val_{\beta}(\bar{t}_0) \mid \begin{array}{l} val_{\beta}(\varphi) = true \ and \ val_{\beta}(\bar{t}_i) \in R_i \ for \ i = \{1, \dots, n\} \\ for \ some \ variable \ assignment \ \beta \end{array} \right\}.$$

$$(4)$$

The algorithm starts with an empty ARG, and then adds nodes and edges to it until the ARG is *complete*, which intuitively means that all possible non-trivial edges are present in the graph. To define the notion of a complete ARG, we first need to characterize what it means for a clause to correspond to a feasible edge of the ARG:

Definition 4. A clause $h \in C$ is feasible for nodes $v_1, \ldots, v_n \in V$ of an ARG (V, E) if

- h is of the form $\forall \bar{x}. \ H \leftarrow p_1(\bar{t}_1) \land \cdots \land p_n(\bar{t}_n) \land \varphi;$
- nodes v_1, \ldots, v_n correspond to the body of h, i.e., for $i \in \{1, \ldots, n\}$ it is the case that $v_i = (p_i, R_i)$ for some $R_i \in \mathcal{R}_{p_i}$;
- the constraints imposed by φ and v_1, \ldots, v_n are not contradictory, i.e., there is a variable assignment β such that $val_{\beta}(\varphi) = true$ and $val_{\beta}(\bar{t}_i) \in R_i$ for $i = \{1, \ldots, n\}$.

Definition 5. An ARG (V, E) is complete for a set C of CHCs if

- for every CHC $h \in C$ that has a head different from false, and that is feasible for $v_1, \ldots, v_n \in V$, there is some edge $\langle (v_1, \ldots, v_n), h, v_0 \rangle \in E$;
- there is no CHC $h \in C$ with head false that is feasible for any $v_1, \ldots, v_n \in V$.

We can finally observe that complete ARGs correspond to models of the clause set \mathcal{C} .

Lemma 1. A set C of CHCs has a complete ARG iff C is satisfiable.

Algorithm 1 describes how ARGs can be constructed for a given clause set \mathcal{C} . The algorithm starts with an empty ARG (V, E), and maintains a queue $Q \subseteq \mathcal{C} \times V^*$ of feasible edges to be added to the graph next. The queue is initialized with the clauses with empty body, representing the initial states of a program. Once the queue runs empty, the constructed ARG is complete and the set \mathcal{C} has been shown to be satisfiable.

In each iteration, in lines 5–6 an element (h, \bar{v}) is picked and removed from the queue Q. If the head of h is *false* (line 7), the edge to be added might be part of a witness of unsatisfiability of C. In this case, it has to be checked whether the nodes \bar{v} are over-approximate (line 8); this can happen when the containment

8

Input: A set C of CHCs $\mathbf{Output}: \ Satisfiability \ of \ \mathcal{C}$ **Initialise**: $V := \emptyset, E := \emptyset, Q := \{(h, ()) \mid h \in \mathcal{C}, h = \forall \bar{x}. H \leftarrow \varphi\}$ 1 while true do $\mathbf{2}$ if Q is empty then return satisfiable 3 else 4 Pick $(h, \bar{v}) \in Q$ to be considered next (guided by GNNs) $\mathbf{5}$ $Q := Q \setminus \{(h, \bar{v})\}$ 6 if the head in h is false then $\mathbf{7}$ if derivation of false is genuine then 8 return unsatisfiable 9 else 10 Refine over-approximations 11 Delete all affected nodes in (V, E)12Regenerate elements in Q13 end 14 15else Assume $h = \forall \bar{x}. \ p_0(\bar{t}_0) \leftarrow p_1(\bar{t}_1) \land \cdots \land p_n(\bar{t}_n) \land \varphi$ $\mathbf{16}$ Compute new node $u = (p_0, R_0)$ for (h, \bar{v}) 17 if $u \notin V$ then 18 $V := V \cup \{u\}$ 19 $Q := Q \cup \left\{ (d, (w_1, \dots, w_m)) \mid \begin{array}{l} d \in \mathcal{C} \text{ is feasible for } w_1, \dots, w_m, \\ u \in \{w_1, \dots, w_m\} \subseteq V \end{array} \right\}$ 20 end 21 $\mathbf{22}$ $E := E \cup \{(\bar{v}, h, u)\}$ end 23 end $\mathbf{24}$ 25 end

Algorithm 1: Abstract symbolic model checking algorithm for checking satisfiability of CHCs

in (4) is sometimes strict in the constructed ARG, and occurs in particular when instantiating the abstract algorithm as CEGAR (see Section 3.2). The over-approximation can then be refined in lines 11–13.

If the head of h is not *false*, a further edge is added to the graph by computing in line 17 some set R_0 satisfying (4). If the resulting node u is new in the graph, the queue Q is updated by adding possible outgoing edges for u (line 19–20).

In this paper, we apply the GNN-based guidance in line 5. Specifically, we presume that the GNN can predict the probability of h being in MUSes for each $q = (h, \bar{v}) \in Q$. We then combine this probability with certain features of q, such as the number of iterations $q \in Q$ has been waiting, to calculate a priority of q. When selecting an element q from Q, we consult this priority. We explain more details in Section 4.

3.2 CEGAR- and Symbolic Execution-Style Exploration

We now discuss two concrete instantiations of Algorithm 1. The first one, in this paper called SymEx, resembles the symbolic execution [27] of a program, and represents the relation R in an ARG node (p, R), for a k-ary relation symbol p, as a formula over free variables z_1, \ldots, z_k .

To obtain SymEx, in line 17 of Algorithm 1 the relation R_0 is derived from the nodes \bar{v} by simple symbolic manipulation. Assuming that $v_i = (p_i, R_i)$ for $i \in \{1, ..., n\}$, and the relations R_i are all represented as formulas, we can define:

$$R_0[z_1, \dots, z_k] = \exists \bar{x}. \ \bar{z} = \bar{t}_0 \land R_1[\bar{t}_1] \land \dots \land R_n[\bar{t}_n] \land \varphi$$

where the notation $R_i[\bar{t}_i]$ expresses that the terms \bar{t}_i are substituted for the free variables \bar{z} . In our implementation on top of the CHC solver Eldarica [28], the formula $R_0[z_1, \ldots, z_k]$ is afterwards simplified by eliminating the quantifiers, albeit only in a best-effort way by running the built-in formula simplifier of Eldarica. In SymEx, since no over-approximation is applied, the test in line 8 always succeeds, and lines 11–13 are never executed.

Our second instantiation, called CEGAR, is designed following counterexampleguided abstraction refinement [5,29] with Cartesian predicate abstraction [30]. In this version of the algorithm, we assume that a finite pre-defined set Π_p of predicates is available for every relation symbol p. If p is k-ary, then the elements of Π_p are formulas over free variables z_1, \ldots, z_k . The relation R in a node (p, R)is now represented as a subset of Π_p .

In line 17, the set R_0 is computed by determining the elements of Π_{p_0} that are entailed by the body of clause h:

$$R_0 = \{ \phi \in \Pi_{p_0} \mid \bar{z} = \bar{t}_0 \land \bigwedge R_1[\bar{t}_1] \land \dots \land \bigwedge R_n[\bar{t}_n] \land \varphi \models \phi \}$$

The notation $\bigwedge R_i[\bar{t}_i]$ denotes the conjunction of the elements of R_i , with \bar{t}_i substituted for the free variables \bar{z} .

Over-approximation in CEGAR stems from the fact that a chosen set of predicates Π_p will oftentimes not be able to exactly represent a relation; the constructed ARG might then include facts $p(\bar{a})$ that are not actually entailed by C. In line 8, the algorithm therefore has to verify that discovered derivations of *false* are genuine. This is done by collecting the clauses that were used to derive the nodes \bar{v} in the ARG and constructing a counterexample tree. If the counterexample turns out to be *spurious*, further predicates are added to the sets Π_p , for instance using tree interpolation [31], in lines 11–13.

4 Guiding CHC Solvers using MUSes

In this section, we begin by defining the notion of Minimal Unsatisfiable Sets (MUSes), then we detail the process of collecting three types of training labels using the MUSes. Following that, we explain the various strategies of employing the predicted probability of a CHC being in MUSes to guide Algorithm 1.

4.1 Minimal Unsatisfiable Sets

Throughout the section, assume that C is an unsatisfiable set of CHCs.

Definition 6 (Minimal Unsatisfiable Set). A subset $U \subseteq C$ is a Minimal Unsatisfiable Set (MUS) if U is unsatisfiable and for all CHCs $h \in U$ it is the case that $U \setminus \{h\}$ is satisfiable.

Intuitively, MUSes of a set of CHCs encoding a program correspond to minimal counterexamples (i.e., a subset of program statements) witnessing the incorrectness of the program. MUSes are therefore good candidates for guiding CHC solvers towards the critical clauses, and we aim at predicting MUSes using GNNs.

The number of MUSes can, however, be exponential in the number of CHCs in C. We therefore consider the *union*, *intersection*, and a particular *single* MUS for C. Denoting the set of all MUSes of C by MUS(C), those are:

$$\mathcal{C}_{\text{MUSes}}^{\text{union}} = \bigcup_{\text{MUS}} \text{MUS}(\mathcal{C}), \qquad \mathcal{C}_{\text{MUSes}}^{\text{intersection}} = \bigcap_{\text{MUS}} \text{MUS}(\mathcal{C}),$$
$$\mathcal{C}_{\text{MUSes}}^{\text{single}} = \underset{U \in \text{MUS}(\mathcal{C})}{\operatorname{argmax}} \text{numAtom}(U),$$

where numAtom(U) is the total number of atoms of the CHCs in U, and C_{MUSes}^{single} is some MUS that maximizes the total number of atoms. The three clause sets can be computed using the OptiRica extension of the Eldarica Horn solver [32].

Intuitively, C_{MUSes}^{union} includes all information about possible MUSes and encourages the algorithm to go through all possible error-prone areas. In contrast, $C_{MUSes}^{intersection}$ only takes the intersection of all MUSes which can guide the algorithm to only focus on the most suspicious areas. C_{MUSes}^{single} is one of MUSes and corresponds to a long path in the ARG, given that a high number of atoms is associated with a large number of nodes. We believe a long path contains intricate information, which is challenging for human to parse, but easier for a deep-learning-based model to find.

We form three types of Boolean clause labelings by using C_{MUSes}^{union} , $C_{MUSes}^{\text{intersection}}$, and $C_{MUSes}^{\text{single}}$, respectively. For instance, for C_{MUSes}^{union} , we obtain the labels

$$l^{\text{union}}(c) = \begin{cases} 1 & \text{if } h \in \mathcal{C}_{\text{MUSes}}^{\text{union}} \\ 0 & \text{if } h \in \mathcal{C} \setminus \mathcal{C}_{\text{MUSes}}^{\text{union}} \end{cases}$$

4.2 MUS-Based Guidance for CHC Solvers

Eldarica Heuristics The implementations of CEGAR and SymEx in the Horn solver Eldarica [28] by default use fixed, hand-written selection heuristics in line 5 of Algorithm 1. Such heuristics are defined by a ranking function $r: Q \to \mathbb{Z}$ that maps every element in the queue to an integer; in line 5, the element of Q is picked that minimizes r. The standard implementation of r used for CEGAR is defined by

$$EldCEGARRank(q) = numPredicate(q) + birthTime(q) + falseClause(q)$$
,

where $numPredicate((h, \bar{v})) = \sum_{(p,R)\in\bar{v}} |R|$ is the total number of predicates occurring in the considered nodes of the ARG; birthTime(q) is the iteration (as an integer) in which q was added to Q;³ and $falseClause((h, \bar{v}))$ is 0 if the head of h is not false, and some big integer otherwise. The rationale behind the ranking function is that nodes with few predicates tend to subsume nodes with many predicates, every clause should be picked eventually, and clauses with head false will trigger either termination of the algorithm or abstraction refinement.

In SymEx, for a node (p, R), we define numConstraint(R) to be the number of conjuncts of R. For instance, for $R = (x > 1 \land y < 0)$ we would get numConstraint(R) = 2. The default Eldarica ranking function for SymEx is defined by

 $EldSymExRank((h, \bar{v})) = \sum_{(p,R)\in\bar{v}} numConstraint(R)$.

Similar to *numPredicate*, the intuition behind *EldSymExRank* is that nodes with larger formulas (more restrictions) tend to be subsumed by nodes with smaller formulas (fewer restrictions).

MUS-guided Heuristics We now introduce several new ranking functions defined with the help of MUSes. For this, suppose that C is the set of CHCs that a Horn solver is applied on. Under the assumption that C is unsatisfiable, we obtain three labeling functions l^{union} , $l^{\text{intersection}}$, l^{single} that are able to point out clauses to be prioritized in Algorithm 1.

In practice, of course, the status of C will initially be unknown. We therefore use three GNNs to *predict* labels $l^{\text{union}}(h)$, $l^{\text{intersection}}(h)$, $l^{\text{single}}(h)$, respectively, given just the set C and some clause $h \in C$ as input. We interpret the prediction of a GNN as a *probability* P(h) of a clause h to be in the union, intersection, or the single MUS set, and use those probabilities to define ranking functions. Table 1 lists several candidate ranking functions in terms of this membership probability P, where P stands for one of P^{union} , P^{single} , or $P^{\text{intersection}}$.

We consider two ways to convert probabilities to integers that can be used in the ranking function. In $rank_P(q)$, the elements $q = (h, \bar{v})$ of the queue Q are first sorted in descending order of P(h); the number $rank_P(q)$ is the position of q in this sequence. This means that $rank_P(q)$ ranges from 0 to |Q| - 1, and elements with large probability P(h) will be assigned small rank.

In $norm_P(q)$, we assume that $\min_P = \min\{P(h) \mid (h, \bar{v}) \in Q\}$ and $\max_P = \max\{P(h) \mid (h, \bar{v}) \in Q\}$ denote the minimum and maximum probability, respectively, among elements of Q. The normalized value $norm_P(q) \in [0, 1]$ is defined by

$$norm_P((h, \bar{v})) = \frac{P(h) - \min_P}{\max_P - \min_P}$$

In the ranking functions, we multiple $norm_P(q)$ with a negative coefficient *coef* and round the result to the nearest integer.

³ Strictly speaking, this information cannot be computed from q, it is in practice stored as an additional field of the queue elements.

Algorithm	Name	Ranking function
	Fixed	EldCEGARRank(q)
	Random	RandomRank
CEGAR	Score	$coef \cdot norm_P(q)$
	Rank	$rank_P(q)$
	R-Plus	$rank_P(q) + EldCEGARRank(q)$
	S-Plus	$coef \cdot norm_P(q) + EldCEGARRank(q)$
	R-Minus	$rank_P(q) - EldCEGARRank(q)$
	S-Minus	$coef \cdot norm_P(q) - EldCEGARRank(q)$
	Fixed	EldSymExRank(q)
	Random	RandomRank
SymEy	Score	$coef \cdot norm_P(q)$
Sym	Rank	$rank_P(q)$
	R-Plus	$rank_P(q) + EldSymExRank(q) + birthTime(q)$
	S-Plus	$coef \cdot norm_P(q) + EldSymExRank(q) + birthTime(q)$
	R-Minus	$rank_P(q) - EldSymExRank(q) - birthTime(q)$
	S-Minus	$coef \cdot norm_P(q) - EldSymExRank(q) - birthTime(q)$
	Two-queue	$\begin{cases} \text{R-Minus,} 80\% \text{ probability} \end{cases}$
	-	[Random, 20% probability

Table 1: Ranking function for queue elements $q \in Q$ in Algorithm 1, where $P \in \{P^{\text{union}}, P^{\text{single}}, P^{\text{intersection}}\}$.

The RandomRank function ensures that each CHC has an equal opportunity to be selected in each iteration. The Two-queue case denotes a setup with two queues, Q_1 and Q_2 , each used with a certain probability in line 5 of Algorithm 1. We list just one such combination, which alternates between queues with the R-Minus and Random functions: there's an 80% chance we use the R-Minus queue and a 20% chance we use the Random queue.

5 Design of Model

As shown in Figure 1, within the CHC-R-HyGNN framework, we first encode a set of CHCs into a graph format, and then we build a GNN model consisting of an encoder, GNN layers, and a decoder.

5.1 Encode CHCs to Graph Representation

We apply the two (hyper-)graph representations of CHCs defined in [1]. We will briefly describe their main features here.

The first graph representation is called a *constraint graph* (CG). This graph encapsulates syntactic information by using nodes to represent each symbol and connecting them with binary edges. Each CHC and each predicate symbol is represented by a unique node. Terms and formulas are represented using their

abstract syntax trees (AST). CHC nodes are connected to the constituent atoms and constraints by binary edges. Within one CHC, common sub-expressions are represented by the same nodes. Different hyper-graph node and edges types are used to distinguish the various encoded operators (see Section 2.3).

The second graph representation is the *control-* and data-flow hypergraph (CDHG). This graph is designed to capture both control- and data-flow within CHCs using hyperedges, and therefore captures the semantics of CHCs more directly than GCs. Similar to the CG, in the CDHG, each CHC is represented by a unique node, and the atoms are rendered in the same way as in CG. Unlike the CG, the CDHG uses *control-flow hyperedges* (CFHEs) to describe the control-flow from the body to the head in each CHC, guarded by the constraint of the CHC. Furthermore, the CDHG uses *data-flow hyperedges* (DFHEs) to represent data-flow from the terms in the body to the terms in the head. These data-flows are also guarded by the constraint.

5.2 Model Structure

The R-HyGNN model consists of three sub-components: (i) encoder, (ii) R-HyGNN [1] (a GNN), and (iii) decoder:

(i) $\mathcal{H}_0 = \operatorname{encoder}(V, \lambda_V, L_V),$ (ii) $\mathcal{H}_t = \operatorname{R-HyGNN}(\mathcal{H}_0, E, L_E),$ (iii) $\hat{\mathcal{L}} = \operatorname{decoder}(\mathcal{H}_t).$

The encoder in (i) first maps each node in V to an integer according to the node's type determined by λ_V and L_V . Then, it passes the encoded integers to a single-layer neural network (embedding layer) to compute initial node representations \mathcal{H}_0 . The R-HyGNN in (ii) is a GNN with its node representation updating rule defined in (3). It takes the initial node representations (\mathcal{H}_0), edges (E), and edge types (L_E) as input and outputs the updated node representations \mathcal{H}_t . The decoder in (iii) first identifies the node that denote the CHC instead of the variables, atoms, or other elements in the CHC, then we collect the representations of these CHC nodes $\mathcal{H}_t^{\text{CHCs}}$ from all node representations \mathcal{H}_t . Finally, we pass $\mathcal{H}_t^{\text{CHCs}}$ to a set of fully connected neural networks to compute the probability of each CHC being in the MUSes, and $\hat{\mathcal{L}}$ is a set of probabilities.

The parameters of neural networks in (i), (ii), and (iii) are optimized together by minimizing the binary cross-entropy loss [33] between $\hat{\mathcal{L}}$ and the true labels \mathcal{L} , i.e.,

$$loss = -\frac{1}{N} \sum_{i=1}^{N} \mathcal{L}_i log(\hat{\mathcal{L}}_i) + (1 - \mathcal{L}_i) log(1 - \hat{\mathcal{L}}_i).$$
(5)

6 Evaluation

We first describe how the benchmarks are split for training and evaluation, then list some important parameters. Finally, we show and explain the experimental results. This work can be reproduced by following the instructions in [34].

Linear LIA problems					Non-linear LIA problems					
8705					8425					
Benchmarks for	ing	Holdo	out set	Benchmarks for	Holde	Holdout set				
7834 (90%)			871 ((10%)	7579 (90%)			846 ((10%)	
UNSAT	SAT	T/O	Eval.	N/A	UNSAT	SAT	T/O	Eval.	N/A	
1585	4004	2245	383	488	3315	4010	254	488	358	
Train Valid N/A					Train Valid N/A					
782 87 716					1617 180 1518]				

Table 2: Distribution of the number of problems for both training and evaluation. T/O, N/A, and Evail. denote timeout, not available, and evaluation, respectively.

6.1 Benchmark

The training and evaluation data are specified in Table 2, and available in [20]. There are 8705 linear and 8425 non-linear LIA problems, taken from CHC-COMP 2021 [19]. We first uniformly reserved 10% of the benchmarks as holdout set for the final evaluation. We ran the CEGAR in Eldarica using a 3-hour timeout to solve the remaining 90% of benchmarks, leading to three groups of benchmarks: SAT, UNSAT, and timeout. For UNSAT problems, we also computed the MUS sets C_{MUSes}^{innion} , $C_{MUSes}^{intersection}$, and C_{MUSes}^{single} . Some problems in UNSAT were eliminated in this process (N/A, for both training and evaluation) because the problems were trivial (already solved by the Eldarica preprocessor), the process of extracting MUSes timed out (3 hours), or a timeout occurred when encoding CHCs as graphs. The remaining problems are divided into training (90%) and validation (10%) datasets.

6.2 Parameters

We select the hyper-parameters for the GNN empirically, and according to the experimental results from [1]. We set the vector size of the initial node representation and the number of neurons in all intermediate neural network layers to 32; we also set the number of message-passing steps to 8 (i.e., applying (3) 8 times). The constant *coef* in Table 1 is -1000. Other parameters and the instructions to reproduce these results can be found in [20].

6.3 Experimental Results

In our experiment, we measure the number of solved problems and the average solving time for the holdout evaluation set. This included 383 and 488 problems in the linear and non-linear LIA datasets, respectively. The timeout for evaluating each problem is 1200 seconds. The additional overhead for reading and applying the GNN predicted results in each iteration is included in the solving

Table 3: Overview of the best ranking function and improvement in number of solved problems compared to the Eldarica. A ranking function marked with * (e.g., S-Plus*) denotes that there are multiple ranking functions with the same performance.

Bonchmark	MUS	Best ranking function (improvement in $\%$)								
Algorithm	data set	Number	of Solved I	Problems	Average Time					
	(best count)	Total	SAT	UNSAT	All	Common	SAT	δ) 2 Γ UNSAT nus Rank %) (31.1%) nus Rank %) (36.3%) nus S-Plus %) (26.5%) om S-Plus %) (17.6%) om R-Plus %) (17.6%) om R-Plus %) (16.8%) re R-Minus %) (16.8%) re R-Plus %) (16.8%) re R-Plus %) (19.9%) nus Two-Q %) (11.2%) -Q S-Plus		
	Union	R-Plus	R-Plus	R-Minus	R-Plus	S-Plus	S-Minus	Rank		
Linear	(0)	(1.4%)	(2.4%)	(1.0%)	(1.3%)	(19.1%)	(46.5%)	(31.1%)		
CEGAR	Single	Rank	R-Plus	Rank	R-Plus	S-Plus	R-Minus	Rank		
	(3)	(3.6%)	(4.0%)	(1.0%) (1. Rank R- (8.2%) (1. R-Plus R- (9.3%) (3. Random Tw (2.0%) (0. Random Ran (2.0%) (0. S-Plus* S- (0.0%) (1. S-Plus* S- (0.0%) (6. S-Plus* R- (0.0%) (5.	(1.9%)	(26.6%)	(57.9%)	(36.3%)		
	Intersection	R-Plus	S-Plus	R-Plus	R-Plus	S-Plus	R-Minus	S-Plus		
	(4)	(4.1%)	(0.8%)	(9.3%)	(3.1%)	(27.6%)	(45.0%)	(0.0%)		
	Union	Two-Q	S-Plus*	Random	Two-Q	R-Minus	R-Minus	S-Plus		
Linear	(4)	(1.0%)	(0.0%)	(2.0%)	(0.9%)	(12.7%)	(30.2%)	(26.5%)		
SymEx	Single	S-Minus*	S-Plus*	Random	Random	S-Plus	Random	S-Plus		
	(3)	(0.5%)	(0.0%)	(2.0%)	(0.8%)	(12.9%)	(28.4%)	(17.6%)		
	Intersection	S-Plus*	S-Plus*	S-Plus*	S-Plus	Score	Random	R-Plus		
	(5)	(1.0%)	(0.0%)	(2.0%)	(1.3%)	(9.5%)	(28.4%)	(35.8%)		
Non	Union	S-Plus	S-Plus	S-Plus*	S-Plus	R-Minus	Rank	S-Plus		
Lincer	(7)	(0.5%)	(0.8%)	(0.0%)	(7.1%)	(20.8%)	(53.5%)	(19.4%)		
	Single	R-Plus	R-Plus	R-Plus*	R-Plus	S-Plus	R-Minus	$\operatorname{R-Minus}$		
CEGAN	(1)	(0.2%)	(0.4%)	(0.0%)	(6.6%)	(18.4%)	(52.8%)	(14.2%)		
	Intersection	R-Plus*	S-Plus	S-Plus*	R-Plus	R-Plus	Rank	S-Plus		
	(1)	(0.0%)	(0.5%)	(0.0%)	(5.9%)	(20.3%)	(45.8%)	(16.8%)		
Non	Union	Two-Q	S-Minus*	Random	Two-Q	R-Minus	Score	R-Plus		
Lincer	(6)	(6.1%)	(1.6%)	(12.3%)	(13.3%)	(7.3%)	(5.1%)	(19.9%)		
SymEy	Single	Two-Q	Score	Two-Q	Two-Q	Rank	R-Minus	Two-Q		
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	(12.4%)	(-2.2%)	(0.2%)	(11.2%)						
	Intersection	$\mathbf{Two-Q}$	S-Plus	Two-Q	Two-Q	S-Minus	Two-Q	S-Plus		
	(3)	(6.1%)	(1.6%)	(12.9%)	(12.7%)	(0.6%)	(1.7%)	(5.4%)		

time. The numerical results are shown in Tables 3 and 4. We also visualize some numerical results by scatter plots in Figure 2.

In Tables 3 and 4, under the Number of Solved Problems column, the Total, and SAT, UNSAT columns denote the number of totals solved, solved SAT, and solved UNSAT problems, respectively. Under the Average Time column, the All column denotes the average solving time for all problems, including those that timed out; the Common column means the average solving time for problems that were commonly solved using one of the ranking functions in Table 1, and the default Eldarica ranking function; the SAT and UNSAT columns are the average solving times for SAT and UNSAT problems, respectively. In certain cells, the percentage in brackets represents the improvement compared to the corresponding default ranking function. The bold text highlights the best performance in a Benchmark Algorithm block for each measurement.

		Numbe	er of Sol	ved Problems	Average Time (improvement %) All Common SAT UNSA 519.38 25.77 38.97 8.77 523.58 27.49 37.05 15.8 (-0.8%) (-29.5%) (4.9%) (-80.77) 512.41 21.65 42.89 11.9 (1.3%) (16.0%) (-10.1%) (-36.7)			
			mprove	ment 70)		(improv	ement 70)
Benchmark Algorithm	Ranking Function	Total	SAT	UNSAT	All	Common	SAT	UNSAT
	Default	222	125	97	519.38	25.77	38.97	8.77
	Random	221	124	97	523.58	27.49	37.05	15.85
Linear	rtandom	(-0.5%)	(-0.8%)	(0.0%)	(-0.8%)	(-29.5%)	(4.9%)	(-80.7%)
CEGAR	R-Plus	225	128	97	512.41	21.65	42.89	11.99
		(1.4%)	(2.4%)	(0.0%)	(1.3%)	(16.0%)	(-10.1%)	(-36.7%)
	R-Minus	220	122	98	526.08	18.02	30.93	21.60
		(-0.9%)	(-2.4%)	(1.0%)	(-1.3%)	(-24.4%)	(20.6%)	(-146.3%)
	S-Plus	222	125	97	517.43	20.92	34.13	7.32
		(0.0%)	(0.0%)	(0.0%)	(0.4%)	(19.1%)	(12.4%)	(16.5%)
	S-Minus	219	122	97	522.97	12.50	20.80	9.81
		(-1.4%)	(-2.4%)	(0.0%)	(-0.770) 502.16	(2.470)	(40.3%)	(-11.9%)
	Portfolio	(2.2%)	(4.0%)	99 (2.1%)	(2.1%)	(20.1%)	43.07	(19.94)
	Default	200	(4.070)	(2.170)	500.68	22.16	(-17.270)	(-127.470)
	Delault	200	101	101	586.12	30.08	30.60	20.05
	Random	(0.5%)	(-1.0%)	(2.0%)	(0.8%)	(-8.5%)	(28.4%)	(-100.7%)
Linear		192	101	91	617.60	38.59	52.87	21.99
SymEx	R-Plus	(-4.0%)	(0.0%)	(-8.1%)	(-4.6%)	(-10.9%)	(4.6%)	(-110.6%)
		200	100	100	586.24	24.67	38.69	10.60
	R-Minus	(0.0%)	(-1.0%)	(1.0%)	(0.8%)	(12.7%)	(30.2%)	(-1.5%)
		198	101	97	595.02	30.22	50.97	7.67
	S-Plus	(-1.0%)	(0.0%)	(-2.0%)	(-0.7%)	(11.6%)	(8.0%)	(26.5%)
	a	201	101	100	586.35	30.64	50.57	10.65
	S-Minus	(0.5%)	(0.0%)	(1.0%)	(0.7%)	(7.8%)	(8.8%)	(-2.0%)
	-	202	101	101	585.58	35.11	49.94	20.14
	Two-queue	(1.0%)	(0.0%)	(2.0%)	(0.9%)	(-5.9%)	(9.9%)	(-92.9%)
	Doutfolio	206	101	105	569.1	25.79	44.58	10.16
	Fortiono	(3%)	(0.0%)	(6.1%)	(3.7%)	(22.2%)	(19.6%)	(2.6%)
	Default	432	250	182	131.12	42.05	43.34	40.28
Non	Bandom	425	243	182	143.42	34.27	34.84	38.75
Linear	manuom	(-1.6%)	(-2.8%)	(0.0%)	(-9.4%)	(-11.1%)	(19.6%)	(3.8%)
CEGAR	B-Plus	432	250	182	122.29	31.74	28.59	37.82
020/11	101140	(0.0%)	(0.0%)	(0.0%)	(6.7%)	(17.8%)	(34.0%)	(6.1%)
	R -Minus	417	240	177	154.07	26.20	21.46	32.51
		(-3.5%)	(-4.0%)	(-2.7%)	(-17.5%)	(20.8%)	(50.5%)	(19.3%)
	S-Plus	434	252	182	121.75	34.04	35.97) UNSAT 8.77 15.85 (-80.7%) 11.99 (-36.7%) 21.60 (-146.3%) 9.81 (-11.9%) 19.94 (-127.4%) 10.44 20.95 (-100.7%) 21.99 (-10.6%) 10.60 (-1.5%) 7.67 (26.5%) 10.65 (-2.0%) 20.14 (-92.9%) 10.65 (-2.0%) 20.14 (-92.9%) 10.65 (-2.0%) 20.14 (-92.9%) 10.16 (2.6%) 33.75 (3.8%) 37.82 (6.1%) 32.51 (19.3%) 39.10 (2.9%) 39.10 (2.9%) 38.95 (3.3%) 31.75 (21.2%) 27.59 41.83 (-51.6%) 27.59 41.83 (-51.6%) 27.59 41.83 (-51.6%) 27.59 41.83 (-51.6%) 27.59 41.83 (-51.6%) 27.59 41.83 (-51.6%) 27.59 41.83 (-51.6%) 27.59 41.83 (-5.16%) 27.59 41.83 (-5.16%) 27.59 41.83 (-1.4%) 27.59 41.83 (-1.4%) 27.59 (-1.4%) 26.75 (-3.75%) 32.51 (-17.8%) 26.75 (-3.0%)
		(0.5%)	(0.8%)	(0.0%)	(7.1%)	(13.1%)	(17.0%)	(2.9%)
	S-Minus	(2507)	(2.907)	(1.607)	(19.02)	(2.0%)	20.33 (20.907)	38.90 (9.907)
		(-2.370)	253	(-1.070)	(-13.770)	28.24	30.57	31.75
	Portfolio	(0.7%)	(1.2%)	(0.0%)	(13.4%)	(20.24)	(20.5%)	(21.2%)
	Default	342	187	155	343.82	28.39	29.05	27.59
	Delault	362	188	174	301.90	32.67	36.24	41.83
Non	Random	(5.8%)	(0.5%)	(12.3%)	(12.2%)	(-15.1%)	(-24.8%)	(-51.6%)
Linear		339	190	149	357.18	27.88	47.71	22.10
Non Linear CEGAR Non Linear SymEx	R-Plus	(-0.9%)	(1.6%)	(-3.9%)	(-3.9%)	(0.3%)	(-64.2%)	(19.9%)
- 5	5.10	361	189	172	299.86	26.35	37.68	27.98
	R-Minus	(5.6%)	(1.1%)	(11.0%)	(12.8%)	(7.3%)	(-29.7%)	(-1.4%)
	C Dl	340	189	151	352.84	29.04	41.41	24.54
	S-Plus	(-0.6%)	(1.1%)	(-2.6%)	(-2.6%)	(-0.3%)	(-42.5%)	(11.1%)
	S Minus	362	190	172	303.65	28.62	44.11	37.95
	5-minus	(5.8%)	(1.6%)	(11.0%)	(11.7%)	(-0.4%)	(-51.8%)	(-37.5%)
	Two-quoue	363	189	174	297.93	30.15	41.14	32.51
	• wo-quede	(6.1%)	(1.1%)	(12.3%)	(13.3%)	(-6.2%)	(-41.6%)	(-17.8%)
	Portfolio	366	191	175	288.85	22.29	42.42	26.75
	1 01 01 01 01 0	(7.0%)	(2.1%)	(12.9%)	(16.0%)	(21.4%)	(-46.0%)	(3.0%)

Table 4:	Evaluation	on	holdout	problems	using	union	dataset.	The	time	unit	is
second.											

Table 3 displays the best ranking function and its improvement over the default Eldarica ranking function in different measurements for various combinations of benchmarks (linear and non-linear LIA), algorithms (CEGAR and SymEx), and MUS datasets (union, single, and intersection of MUSes). In the MUS dataset column, the numbers in brackets represent the count of bold text cells in the row, indicating the number of best performances achieved by that type of MUS dataset. For instance, in the last row (i.e., the Intersection row of the Non-Linear SymEx block), the number in the bracket is counted from the bold text highlighted cells in columns Total, SAT, and UNSAT under the Number of Solved Problems. This suggests that using the intersection of the MUSes dataset achieves the best performance when the evaluation set is non-linear and the algorithm is SymEx. Across the entire table, there are 17, 10, and 13 bold text counts for the union, single, and intersection MUS datasets, respectively. This indicates that the union is the most effective MUS dataset for better performance across different benchmarks and algorithms. Consequently, we provide further numerical details for the union MUS dataset in Table 4. Evaluation results for both the single and intersection MUS datasets can be found in [20].

Table 4 illustrates the evaluation results using MUS-guided ranking functions (refer to Table 1), compared to the default and random ranking functions. The portfolio rows in each Benchmark Algorithm block represent the best potential of ranking functions, as they accumulate the best performance for each problem across all ranking functions. The model is trained using the union of MUSes.

In terms of the total number of solved problems, the improvement for the Linear dataset is at most 1.4%, achieved by the CEGAR algorithm with the R-Plus ranking function. Meanwhile, for the Non-linear dataset, the improvement is 6.1%, achieved by the SymEx algorithm with the two-queue ranking function. And, this is consistent with the average solving time for all benchmarks.

The best performances in the Common column are inconsistent with the All column under Average Time. This suggests that the number of newly solved problems has a greater impact on the improvement in the average solving time for all problems than the commonly solved problems.

The improvements in the average solving time for SAT and UNSAT problems are 50.5% and 26.5%, achieved by the CEGAR with R-Minus and the SymEx with S-Plus ranking function, respectively. When combined with the corresponding numbers of total solved problems (i.e., -3.5% and -1.0%), it suggests that these ranking functions either solve the problems quickly or not at all.

Furthermore, Figure 2 shows the solving time scatter plots for the problems from the best configurations in each Benchmark Algorithm block in Table 4. Notably, a majority of the dots lie below the diagonal lines in each scatter plot, indicating the solving time is improved by the MUS-guided ranking function for more than half of the problems. This is consistent with the numerical results.

In summary, for both algorithms in different datasets, there is always at least one of the MUS-guided ranking functions that achieves the best result in terms of all aspects of measurements. Using the predicted probabilities alone (i.e., using the Rank and Score ranking function) performs weaker than other



Fig. 2: All benchmark average solving time scatter plots for best ranking functions in different dataset and algorithms. "above/under" means the number of dots above and under the diagonal line.

MUS-guided ranking functions that combine the predicted probability and the default heuristics. Currently, the MUS-guided ranking functions in Table 1 are designed by simply varying the relation symbols "+" and "-" between different elements (e.g., ranking functions S-Plus and S-Minus for both CEGAR and SymEx) or by setting the restart point randomly (i.e., ranking function Twoqueue in SymEx). We believe MUSHyperNet has more potential if the ranking functions are designed carefully or learned from some good tasks.

7 Related Work

Machine learning techniques have been adapted in various ways to assist in formal verification. For example, the study in [35] employs Support Vector Machines (SVM) [36] and Gaussian processes[37] to select heuristics for theorem proving. Similarly, [38] introduces the use of a Recurrent Neural Network Based Language Model (RNNLM) to derive finite-state automaton-based specifications from execution traces. In the domain of selecting algorithms for program verification, [39] apply the Transformer architecture [40], while [41] uses kernel-based methods [42]. With the thriving of deep learning techniques, an increasing number of works are utilizing GNNs to learn the features from programs and logic formu-

lae. This trend is attributed to the inherent structure of these languages, which can be naturally represented as graphs and subsequently learned by GNNs. For instance, studies like [14,15], [16,43], and [44] use GNNs [24,45,46] to learn features from graph-represented logic formulas and programs, aiding in tasks such as theorem proving, SAT solving, and loop invariant reasoning.

One closed idea is NeuroSAT [17,16], which trains a GNN to predict the probability of variables appearing in unsat cores. This prediction can guide the variable branching decisions for Conflict-Driven Clause Learning (CDCL) [47]-based SAT solvers. In a similar vein, our study trains a GNN to predict the probability of a CHC appearing in MUSes. This aids in determining the processing sequence of CHCs in Algorithm 1 used for solving a set of CHCs.

8 Conclusion and Future Works

In this study, we train a GNN model to predict the probability of each CHC in a set of CHCs being in the MUSes. We then utilize these predicted probabilities to guide the abstract symbolic model-checking algorithm in selecting a CHC during each iteration. Extensive experiments demonstrate improvements in both the number of solved problems and average solving time when using the MUSguided ranking functions, compared to the default ranking function. This was observed in two instances of the abstract symbolic model checking algorithm: CEGAR and SymEx. We believe that this approach can be extended to other algorithms, as many could benefit from understanding more about the MUSes of a set of CHCs.

There are several ways to further enhance the performance of MUSHyperNet. One of our future work is to integrate the work of manually designing the ranking functions in Table 1 to the learning process. Regarding the GNN model, we believe that incorporating an attention mechanism could bolster its performance, subsequently refining the quality of the predicted probabilities. Another avenue to explore involves integrating the GNN with the solver in a more interactive manner. Instead of predicting something at once and then using them in each iteration, we could query the GNN model to predict something in real-time based on the current context during each iteration.

Acknowledgement. We thank Zafer Esen for providing assistance in using the symbolic execution engine of Eldarica, and Marc Brockschmidt for various discussions on this work. The computations and data handling were enabled by resources provided by the Swedish National Infrastructure for Computing (SNIC) at UPPMAX and C3SE. The research was partially funded by the Swedish Research Council through grant agreement no. 2018-05973, by a Microsoft Research PhD grant, the Swedish Foundation for Strategic Research (SSF) under the project WebSec (Ref. RIT17-0011), and the Wallenberg project UPDATE.

References

- 1. Chencheng Liang, Philipp Rümmer, and Marc Brockschmidt. Exploring representation of Horn clauses using GNNs (extended technique report). *CoRR*, abs/2206.06986, 2022.
- Alfred Horn. On sentences which are true of direct unions of algebras. The Journal of Symbolic Logic, 16(1):14–21, 1951.
- Emanuele De Angelis, Fabio Fioravanti, John P. Gallagher, Manuel V. Hermenegildo, Alberto Pettorossi, and Maurizio Proietti. Analysis and transformation of constrained Horn clauses for program verification. *CoRR*, abs/2108.00739, 2021.
- Nikolaj Bjørner, Arie Gurfinkel, Ken McMillan, and Andrey Rybalchenko. Horn Clause Solvers for Program Verification, pages 24–51. Springer International Publishing, Cham, 2015.
- Edmund Clarke, Orna Grumberg, Somesh Jha, Yuan Lu, and Helmut Veith. Counterexample-guided abstraction refinement. In E. Allen Emerson and Aravinda Prasad Sistla, editors, *Computer Aided Verification*, pages 154–169, Berlin, Heidelberg, 2000. Springer Berlin Heidelberg.
- Aaron R. Bradley. Sat-based model checking without unrolling. In Ranjit Jhala and David Schmidt, editors, Verification, Model Checking, and Abstract Interpretation, pages 70–87, Berlin, Heidelberg, 2011. Springer Berlin Heidelberg.
- Mark H. Liffiton and Karem A. Sakallah. Algorithms for computing minimal unsatisfiable subsets of constraints. *Journal of Automated Reasoning*, 40(1):1–33, Jan 2008.
- Arie Gurfinkel. Program verification with constrained Horn clauses (invited paper). In Sharon Shoham and Yakir Vizel, editors, *Computer Aided Verification*, pages 19–29, Cham, 2022. Springer International Publishing.
- Karel Chvalovský, Jan Jakubův, Martin Suda, and Josef Urban. ENIGMA-NG: Efficient neural and gradient-boosted inference guidance for E. In Pascal Fontaine, editor, Automated Deduction – CADE 27, pages 197–215, Cham, 2019. Springer International Publishing.
- Jan Jakubův and Josef Urban. ENIGMA: Efficient learning-based inference guiding machine. In Herman Geuvers, Matthew England, Osman Hasan, Florian Rabe, and Olaf Teschke, editors, *Intelligent Computer Mathematics*, pages 292–302, Cham, 2017. Springer International Publishing.
- Jan Jakubův, Karel Chvalovský, Miroslav Olšák, Bartosz Piotrowski, Martin Suda, and Josef Urban. ENIGMA anonymous: Symbol-independent inference guidingmachine (system description). In Nicolas Peltier and Viorica Sofronie-Stokkermans, editors, *Automated Reasoning*, pages 448–463, Cham, 2020. Springer International Publishing.
- Martin Suda. Vampire with a brain is a good ITP hammer. In Boris Konev and Giles Reger, editors, *Frontiers of Combining Systems*, pages 192–209, Cham, 2021. Springer International Publishing.
- 13. Peter W. Battaglia, Jessica B. Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinícius Flores Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, Çaglar Gülçehre, H. Francis Song, Andrew J. Ballard, Justin Gilmer, George E. Dahl, Ashish Vaswani, Kelsey R. Allen, Charles Nash, Victoria Langston, Chris Dyer, Nicolas Heess, Daan Wierstra, Pushmeet Kohli, Matthew Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu. Relational inductive biases, deep learning, and graph networks. *CoRR*, abs/1806.01261, 2018.

- Mingzhe Wang, Yihe Tang, Jian Wang, and Jia Deng. Premise selection for theorem proving by deep graph embedding. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 2786–2796. Curran Associates, Inc., 2017.
- Aditya Paliwal, Sarah M. Loos, Markus N. Rabe, Kshitij Bansal, and Christian Szegedy. Graph representations for higher-order logic and theorem proving. *CoRR*, abs/1905.10006, 2019.
- Daniel Selsam and Nikolaj Bjørner. Neurocore: Guiding high-performance SAT solvers with unsat-core predictions. CoRR, abs/1903.04671, 2019.
- Daniel Selsam, Matthew Lamm, Benedikt Bünz, Percy Liang, Leonardo de Moura, and David L. Dill. Learning a SAT solver from single-bit supervision. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019.
- 18. CHC-COMP benchmarks. Accessed: 2023-09-07, https://chc-comp.github.io/.
- Grigory Fedyukovich and Philipp Rümmer. Competition report: CHC-COMP-21. In Hossein Hojjat and Bishoksan Kafle, editors, *Proceedings 8th Workshop on Horn Clauses for Verification and Synthesis, HCVS@ETAPS 2021, Virtual, 28th March 2021*, volume 344 of *EPTCS*, pages 91–108, 2021.
- 20. Repository for the training and evaluation dataset in this work. Accessed: 2023-09-07, https://github.com/ChenchengLiang/Horn-graph-dataset.
- 21. John Harrison. *Handbook of Practical Logic and Automated Reasoning*. Cambridge University Press, 2009.
- 22. Clark Barrett, Pascal Fontaine, and Cesare Tinelli. The SMT-LIB Standard: Version 2.6. Technical report, Department of Computer Science, The University of Iowa, 2017. Available at www.SMT-LIB.org.
- Ian J. Goodfellow, Yoshua Bengio, and Aaron Courville. Deep Learning. MIT Press, Cambridge, MA, USA, 2016. http://www.deeplearningbook.org.
- Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neural message passing for quantum chemistry. *CoRR*, abs/1704.01212, 2017.
- 25. Michael Schlichtkrull, Thomas N. Kipf, Peter Bloem, Rianne vanden Berg, Ivan Titov, and Max Welling. Modeling relational data with graph convolutional networks. In Aldo Gangemi, Roberto Navigli, Maria-Esther Vidal, Pascal Hitzler, Raphaël Troncy, Laura Hollink, Anna Tordai, and Mehwish Alam, editors, *The Semantic Web*, pages 593–607, Cham, 2018. Springer International Publishing.
- 26. Abien Fred Agarap. Deep Learning using Rectified Linear Units (ReLU). arXiv e-prints, page arXiv:1803.08375, March 2018.
- James C. King. Symbolic execution and program testing. Commun. ACM, 19(7):385–394, july 1976.
- Hossein Hojjat and Philipp Rümmer. The ELDARICA Horn solver. In 2018 Formal Methods in Computer Aided Design (FMCAD), pages 1–7, 2018.
- Susanne Graf and Hassen Saidi. Construction of abstract state graphs with PVS. In Orna Grumberg, editor, *Computer Aided Verification*, pages 72–83, Berlin, Heidelberg, 1997. Springer Berlin Heidelberg.
- 30. Thomas Ball, Andreas Podelski, and Sriram K. Rajamani. Boolean and cartesian abstraction for model checking C programs. In Tiziana Margaria and Wang Yi, editors, Tools and Algorithms for the Construction and Analysis of Systems, 7th International Conference, TACAS 2001 Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2001 Genova, Italy, April

22

2-6, 2001, Proceedings, volume 2031 of Lecture Notes in Computer Science, pages 268–283. Springer, 2001.

- Philipp Rümmer, Hossein Hojjat, and Viktor Kuncak. Disjunctive interpolants for Horn-Clause verification. In Natasha Sharygina and Helmut Veith, editors, *Computer Aided Verification - 25th International Conference, CAV 2013, Saint Petersburg, Russia, July 13-19, 2013. Proceedings*, volume 8044 of *Lecture Notes in Computer Science*, pages 347–363. Springer, 2013.
- 32. Hossein Hojjat and Philipp Rümmer. OptiRica: Towards an efficient optimizing Horn solver. In Geoffrey William Hamilton, Temesghen Kahsai, and Maurizio Proietti, editors, Proceedings 9th Workshop on Horn Clauses for Verification and Synthesis and 10th International Workshop on Verification and Program Transformation, HCVS/VPT@ETAPS 2022, and 10th International Workshop on Verification and Program TransformationMunich, Germany, 3rd April 2022, volume 373 of EPTCS, pages 35–43, 2022.
- David R. Cox. The regression analysis of binary sequences. Journal of the Royal Statistical Society. Series B (Methodological), 20(2):215-242, 1958.
- Code repository for reproduce this work. Accessed: 2023-09-07, https://github.com/ChenchengLiang/Relational-Hypergraph-Neural-Network-PyG.
- 35. James Bridge, Sean Holden, and Lawrence Paulson. Machine learning for firstorder theorem proving. *Journal of Automated Reasoning*, 53, 08 2014.
- Corinna Cortes and Vladimir Vapnik. Support-vector networks. Machine learning, 20(3):273–297, 1995.
- Carl Edward Rasmussen and Christopher K. I. Williams. Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning). The MIT Press, 2005.
- Tien-Duy B. Le and David Lo. Deep specification mining. In Proceedings of the 27th ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2018, pages 106–117, New York, NY, USA, 2018. Association for Computing Machinery.
- Cedric Richter and Heike Wehrheim. Attend and represent: A novel view on algorithm selection for software verification. In 2020 35th IEEE/ACM International Conference on Automated Software Engineering (ASE), pages 1016–1028, 2020.
- 40. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, L ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017.
- Cedric Richter, Eyke Hüllermeier, Marie-Christine Jakobs, and Heike Wehrheim. Algorithm selection for software validation based on graph kernels. *Automated Software Engineering*, 27(1):153–186, June 2020.
- Bernhard Scholkopf and Alexander J. Smola. Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond. MIT Press, Cambridge, MA, USA, 2001.
- Daniel Selsam, Matthew Lamm, Benedikt Bünz, Percy Liang, Leonardo de Moura, and David L. Dill. Learning a SAT solver from single-bit supervision. *CoRR*, abs/1802.03685, 2018.
- Xujie Si, Hanjun Dai, Mukund Raghothaman, Mayur Naik, and Le Song. Learning loop invariants for program verification. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Infor*mation Processing Systems, volume 31. Curran Associates, Inc., 2018.

- Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2009.
- 46. Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard S. Zemel. Gated graph sequence neural networks. In Yoshua Bengio and Yann LeCun, editors, 4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings, 2016.
- Joao P. Marques-Silva and Karem A. Sakallah. GRASP: a search algorithm for propositional satisfiability. *IEEE Transactions on Computers*, 48(5):506–521, 1999.

24