



ARENA: Enhancing Abstract Refinement for Neural Network Verification

Yuyi Zhong^(✉), Quang-Trung Ta, and Siau-Cheng Khoo

School of Computing, National University of Singapore, Singapore, Singapore
{yuyizhong, taqt, khoosc}@comp.nus.edu.sg

Abstract. As neural networks have taken on a critical role in real-world applications, formal verification is earnestly needed to guarantee the safety properties of the networks. However, it remains challenging to balance the trade-off between precision and efficiency in abstract interpretation based verification methods. In this paper, we propose an abstract refinement process that leverages the convex hull techniques to improve the analysis efficiency. Specifically, we introduce the double description method in the convex polytope domain to detect and eliminate multiple *spurious* adversarial labels simultaneously. We also combine the new activation relaxation technique with the iterative abstract refinement method to compensate for the precision loss during abstract interpretation. We have implemented our proposal into a verification framework named ARENA, and assessed its effectiveness by conducting a series of experiments. These experiments show that ARENA yields significantly better verification precision compared to the existing abstract-refinement-based tool DeepSRGR. It also identifies falsification by detecting adversarial examples, with reasonable execution efficiency. Lastly, it verifies more images than the state-of-the-art verifier PRIMA.

Keywords: Abstract refinement · Double description method · Neural network verification

1 Introduction

As neural networks have been proverbially applied to safety-critical systems, formal guarantee about the safety properties of the networks, such as robustness, fairness, etc., is earnestly needed. For example, researchers have been working on robustness verification of neural networks, to ascertain that the network classification result can remain the same when the input image is perturbed subtly and imperceptibly during adversarial attacks [1, 2].

There exists sound and complete verification techniques where the robustness property can be ascertained but regrettably at high complexity and execution cost [3, 4]. For better scalability, several incomplete verifiers have been proposed to analyze larger networks with abstract interpretation technique while bearing

The original version of the chapter has been revised. The acknowledgment section have been corrected. A correction to this chapter can be found at

https://doi.org/10.1007/978-3-031-24950-1_18

© The Author(s), under exclusive license to Springer Nature Switzerland AG 2023, corrected publication 2024
C. Dragoi et al. (Eds.): VMCAI 2023, LNCS 13881, pp. 366–388, 2023.

https://doi.org/10.1007/978-3-031-24950-1_17

exactness sacrifices [5–7]. To mitigate this shortcoming, there have been investigations into better convex relaxation [8, 9] or iterative abstract refinement [10] to make up with the precision loss in abstract interpretation techniques.

This work is inspired by the counterexample guided abstraction refinement (CEGAR) method [11] in program analysis, aiming to improve the precision of abstract interpretation results by identifying *spurious counterexamples*: these are examples which appear to have violated desired analysis outcome – due to over-approximated calculation inherent in the abstract interpretation computation – but can be shown to be fake by the refinement method. Proof of the existence of spurious counterexamples can diminish the range of inconclusive results produced by abstract interpretation. In the context of neural network verification, such spurious counterexamples can be conceptualized as adversarial *regions* that are perceived to have lent support (spuriously) on certain adversarial labels; i.e., labels which differ from the designated label in the robustness test.

An existing work that has successfully employed abstract refinement technique to improve the precision of the abstract-interpretation based verification tool is DeepSRGR [10]. That work repetitively selects an adversarial label and attempts to eliminate the corresponding spurious region progressively through iteration of refinements. Technically, it encodes a spurious region as a linear inequality, adds it to the constraint encoding of the network, and employs linear programming with the objective set to optimize the concrete bounds of selected ReLU neurons in the network. This process is repeated until either the spurious region is found to be inconsistent with the encoded network, or time out.

In this paper, we enhance the existing effectiveness of DeepSRGR by introducing *convex hull* techniques (i.e., techniques that observe and conform to convex property) to abstract refinement computation. Together, these techniques facilitate *simultaneous* elimination of multiple spurious regions, and capitalize the *dependencies* among ReLU neurons. Specifically, we tighten the looseness of ReLU approximation during abstract refinement process through a *multi-ReLU convex abstraction technique* (cf. [9]) that captures dependencies within a set of ReLU neurons. Moreover, we leverage a double-description method (cf. [12]) used in convex polytope computation to eliminate multiple spurious regions simultaneously; this circumvents the challenges faced with the application of linear programming technique to optimize disjunction of linear inequalities.

We have implemented our proposed techniques in a CPU-based prototypical analyzer named ARENA (AbstrAct Refinement Enhancer for Neural network verifAtion). In addition to verifying the robustness property of a network with respect to an image, ARENA is also capable of detecting adversarial examples that ascertain the falsification of the network property. We conducted experiments to assess the effectiveness of ARENA against the state-of-the-art tools, including the CPU-based verifiers DeepSRGR [10] and PRIMA [9], and the GPU-based verifier α, β -CROWN [13]. The results show conclusively that ARENA returns an average of 15.8% more conclusive images compared with DeepSRGR while terminates in comparable amount of time; and it also outperforms PRIMA by returning 16.6% more conclusive images. Furthermore, ARENA can verify or

falsify 79.3% images of that of the state-of-the-art complete tool α, β -CROWN on average for selected networks.

We summarize our contributions below:

- ◊ We adapt the double description method proposed in the convex polytope domain [12] to solve disjuncts of constraints in Linear Programming (LP) encoding, allowing us to prune multiple adversarial labels together to increase overall efficiency.
- ◊ We leverage the multi-ReLU convex abstraction in PRIMA [9] to further refine the abstraction in the analysis process to increase verification precision.
- ◊ We utilize the solutions returned by the LP solver to detect adversarial examples and assert property violation when counter-examples are discovered.
- ◊ We conducted experiments comparing our prototypical analyzer ARENA against state-of-the-art verification tools, and demonstrate high effectiveness in our verification framework. To the best of our knowledge, ARENA outperforms the current state-of-the-art approximated methods that run on CPU.

In the remaining part of the paper, we give an illustrative example showing the overall process of our method in Sect. 2, followed by a formal description of our methodologies in Sect. 3. We demonstrate our evaluation process and experimental results in Sect. 4. Section 5 discusses the current limitation, plan for future work and the generalization of our work. We give a literature review in Sect. 6, which contains closely related works with respect to our research scope. Finally, we summarize our work and conclude in Sect. 7.

2 Overview

In this section, we first describe the abstract refinement technique implemented in DeepSRGR [10]. Then, we discuss its limitations and introduce our approach to overcome them. Table 1 displays the notations we use throughout this section.

Table 1. Notations and descriptions of Sect. 2

| | |
|---------------|-------------------------------------|
| Π | The network constraint set/encoding |
| \mathcal{Y} | A potential adversarial region |
| P | The over-approximate convex hull |

2.1 Spurious Region Guided Refinement

DeepSRGR is a sound but incomplete verification method that relies on the polyhedral abstract domain in DeepPoly [7], where the abstract value of each network neuron x_i is designed to contain four elements (l_i, u_i, l_i^s, u_i^s) . The *concrete* lower bound l_i and upper bound u_i pair forms a closed interval $[l_i, u_i]$ that over-approximates all the values that neuron x_i could take. The *symbolic* constraints l_i^s, u_i^s are linear expressions of x_i defined over preceding neurons with the requirement that $l_i^s \leq x_i \leq u_i^s$.

In the following, we illustrate the verification process where the abstract domain is used to verify the robustness property of a fully-connected network with ReLU activation (Fig. 1) w.r.t the input space $I = [-1, 1] \times [-1, 1]$ of 2 input neurons x_1, x_2 . This network has 3 output neurons y_1, y_2, y_3 , corresponding to the three labels L_1, L_2, L_3 that an input in I can be classified as. Here, the robustness property which we aim to verify is that the neural network can always classify the entire input space I as label L_1 , which corresponds to the output neuron y_1 . More specifically, the verifier should be able to prove that the conditions $y_1 - y_2 > 0$ and $y_1 - y_3 > 0$ always hold for the entire input space I .

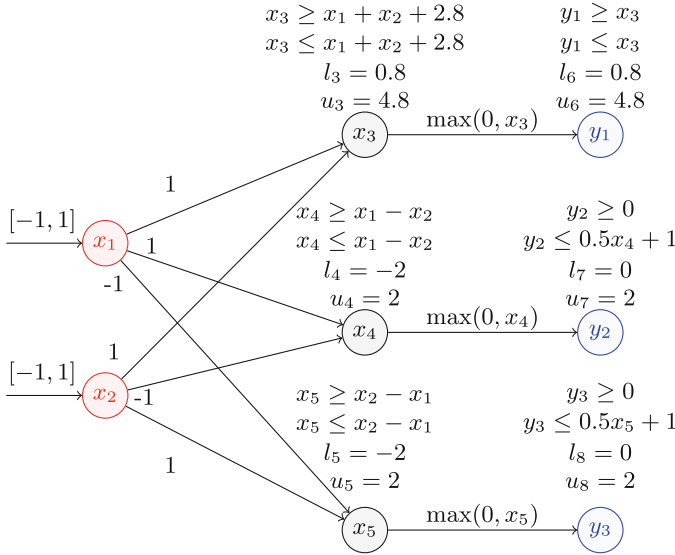


Fig. 1. The example network to perform DeepPoly abstract interpretation

Through the abstract interpretation technique, as deployed by DeepPoly, we can compute the abstract values for each neuron; these are displayed near the corresponding nodes. Specifically, the computed value for the lower bound of $y_1 - y_2$ and $y_1 - y_3$ are both -0.2 (the process of the lower bound computation is provided in Appendix A in our technical report), which fails to assert that $y_1 - y_2 > 0$ and $y_1 - y_3 > 0$. In other word, DeepPoly cannot ascertain the robustness of the network for the given initial input space I . Given the over-approximation nature of abstract interpretation technique, it is not clear if the robustness property can be verified.

In order to further improve the robustness verification of the considered neural network, DeepSRGR conducts a *spurious region guided refinement* process that includes the following steps:

1. Obtain the conjunction of all linear inequities that encode the network, including the input constraint $x_1, x_2 \in [-1, 1]$ and the constraints within the abstract values of all neurons (i.e., the constraints in Fig. 1). We denote this network encoding constraint set as Π .

2. Take the conjunction of the current network encoding and the negation of the property to solve a potential *spurious* region. For example, the feasible region of $\Pi \wedge (y_1 - y_2 \leq 0)$ refers to a *potential* adversarial region (denote as \mathcal{Y}) that may contain a counterexample with adversarial label L_2 (corresponding to output neural y_2); whereas the region *outside* of \mathcal{Y} is already a safe region that will not be wrongly classified as L_2 . However, the region \mathcal{Y} may exist only due to the over-approximate abstraction but does not contain any true counterexample. Therefore, this region is spuriously constructed and could be eliminated. If we successfully eliminate \mathcal{Y} , then we can conclude that label L_2 will not be a valid adversarial label since y_2 never dominates over y_1 .
3. To eliminate the region \mathcal{Y} , DeepSRGR uses the constraints of the region to refine the abstraction using linear programming (LP). For instance, we take Π and $y_1 - y_2 \leq 0$ as the constraint set of linear programming. To obtain tighter bounds for input neurons and unstable ReLU neurons¹, we set the objective function of LP as $\min(x_i)$ and $\max(x_i)$ where $i \in [1, 2, 4]$. The new solved intervals are highlighted in red in Fig. 2, where all the current neuron intervals now specify the region \mathcal{Y} .

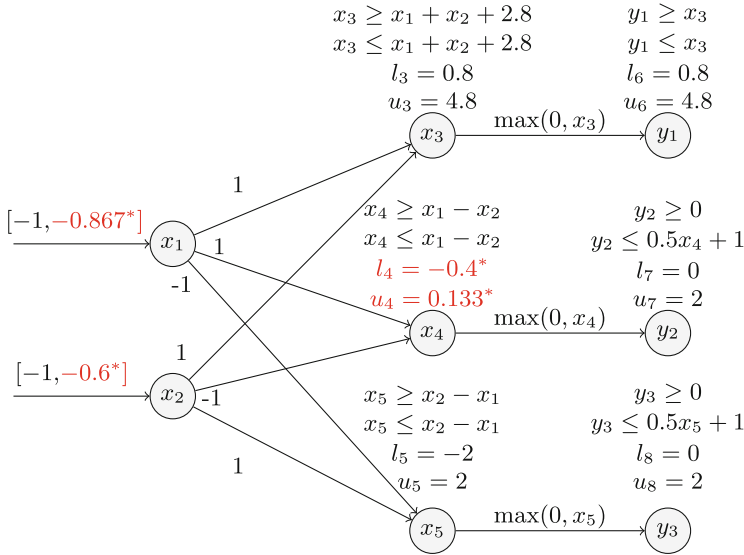


Fig. 2. The effect of applying LP-based interval refinement (in red marked by $*$) (Color figure online)

4. DeepSRGR leverages those tighter bounds to guide the abstract interpretation of the region \mathcal{Y} in the next iteration. It performs a second run on DeepPoly and makes sure that this second run compulsorily follows the new

¹ Unstable ReLU neuron refers to a ReLU neuron whose input range can be both negative and positive (like y_2, y_3).

bounds computed in the previous step. As shown in Fig. 3, the blue colored part refers to the updated abstract values during the second execution of DeepPoly, where the abstraction of all neurons are refined due to the tighter bounds (red colored part) returned by LP solving. Now the lower bound of $y_1 - y_2$ is 0.7, making $y_1 - y_2 \leq 0$ actually infeasible within the region \mathcal{Y} . Therefore, we conclude that \mathcal{Y} is a spurious region that does not contain any true counterexample, and we can eliminate adversarial label L_2 .

5. If we fail to detect $y_1 - y_2 \leq 0$ to be infeasible, DeepSRGR iterates the process from step 2-4 where it calls LP solving and re-executes DeepPoly on the new bounds until it achieves one of the termination conditions: (i) It reaches the maximum number of iterations (DeepSRGR sets it to be 5 by default); or (ii) it detects infeasibility for the spurious region.

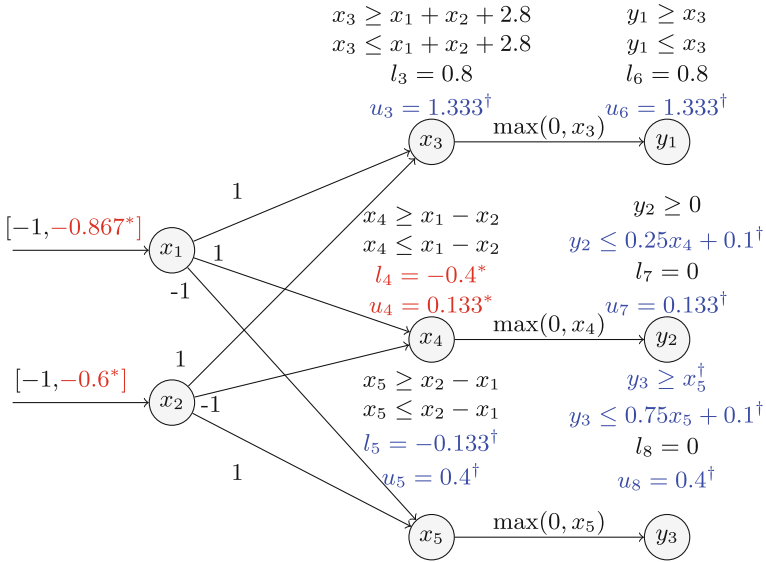


Fig. 3. Results of the second run of DeepPoly (in blue marked by †) (Color figure online)

Similarly, after eliminating the adversarial label L_2 , DeepSRGR will apply the same process to eliminate the spurious region defined by $\Pi \wedge (y_1 - y_3 \leq 0)$, which corresponds to the output neural y_3 and the adversarial label L_3 .

In summary, DeepSRGR uses iterative LP solving and DeepPoly execution to attempt to eliminate spurious regions which do not contain counterexamples. Assuming the ground-truth label to be L_c , DeepSRGR runs this refinement process for each region $\Pi \wedge (y_c - y_t \leq 0)$ where $t \neq c$. If DeepSRGR is able to eliminate all adversarial labels related to output neurons y_t where $t \neq c$, then it successfully ascertains the robustness property of the image. If DeepSRGR fails to eliminate one of the adversarial labels within the iteration boundary, the robustness result remains inconclusive.

2.2 Scaling up with Multiple Adversarial Label Elimination

We mention three contributions in Sect. 1 including efficiency improvement, precision improvement and adversarial example detection. In this section, we only give an overview of our multiple adversarial label elimination method which aims to improve the analysis efficiency; we defer the discussion of the remaining part of our system to Sect. 3.

As mentioned in Sect. 2.1, DeepSRGR invokes the refinement process to sequentially eliminate each spurious region $\Pi \wedge (y_c - y_t \leq 0)$, which corresponds to the adversarial label L_t ($t \neq c$). For an n -label network, it requires $n - 1$ refinement invocations in the worst case, with each invocation taking possibly several iterations. To speed up the analysis, we eliminate multiple spurious regions at the same time in one refinement process.

For example, we aim to detect infeasibility in $\Pi \wedge ((y_1 - y_2 \leq 0) \vee (y_1 - y_3 \leq 0))$ so as to eliminate both adversarial labels L_2 and L_3 simultaneously. The technical challenge behind this *multiple adversarial label* elimination is that linear programming does not naturally support the *disjunction* of linear inequalities. To address this challenge, we compute the over-approximate convex hull P of $(y_1 - y_2 \leq 0) \vee (y_1 - y_3 \leq 0)$ under network encoding Π . As P will be represented as a set of linear inequalities, linear programming is amenable to handle $\Pi \wedge P$.

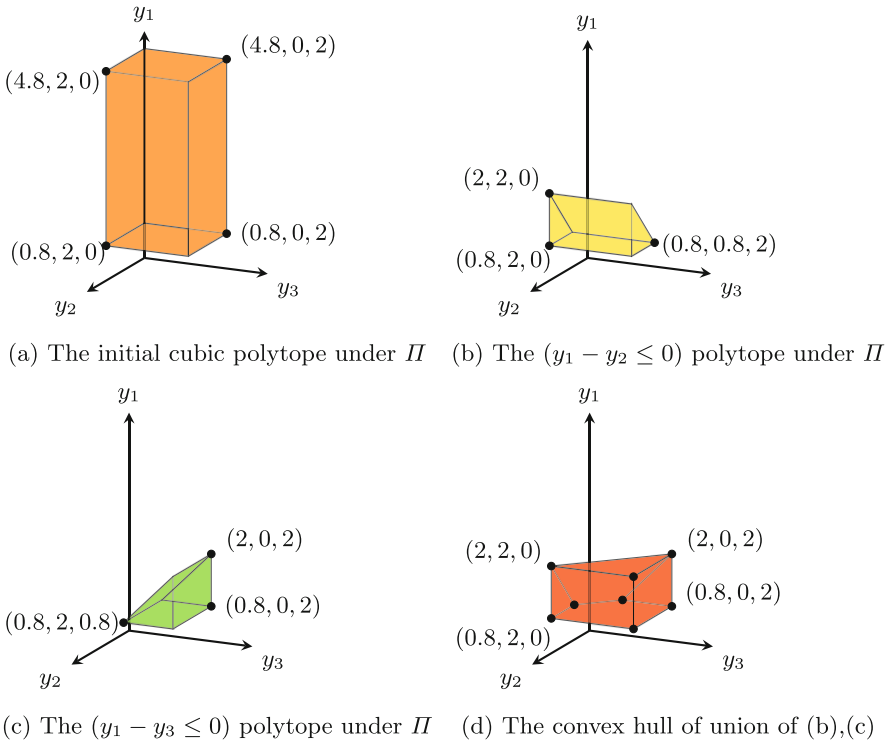


Fig. 4. The convex polytopes under network encoding, with respect to (y_1, y_2, y_3) .

In detail, the initial convex polytope associated with y_1, y_2, y_3 is a 3-D cube pictured in Fig. 4a, where $y_1 \in [0.8, 4.8], y_2 \in [0, 2], y_3 \in [0, 2]$ after we perform DeepPoly as shown in Fig. 1. The convex polytope for the constraint $y_1 - y_2 \leq 0$ under the network encoding Π corresponds to the shape in Fig. 4b where $y_1 - y_2 \leq 0$ is a cutting-plane imposed on the initial cube in Fig. 4a. Similarly, the projection of $y_1 - y_3 \leq 0$ to a convex polytope can be visualized in Fig. 4c. We further compute the over-approximate convex hull P of the union of the two polytopes as in Fig. 4d.

We can observe that P is defined by 8 vertices (annotated as eight black extreme points). It is worth-noting that these 8 vertices actually come from either vertices in Fig. 4b or vertices in Fig. 4c. We will provide the explanation and the theory on how to compute the convex hull of the union of two polytopes in Sect. 3.2. Explicitly, P can also be represented by the following constraint set (1), which correspond to the 7 red-colored surfaces in Fig. 4d:

$$\begin{aligned} -y_1 + y_2 + y_3 \geq 0 \quad y_2 \geq 0 \quad y_3 \geq 0 \quad -1 + 1.25y_1 \geq 0 \\ 2 - y_1 \geq 0 \quad 2 - y_2 \geq 0 \quad 2 - y_3 \geq 0 \end{aligned} \tag{1}$$

We take the network encoding Π and constraint set of P as the input to the LP solver, and conduct interval solving as in Sect. 2.1. We annotate the new bounds obtained through LP solving as red color, and the updated abstract values after the second abstract interpretation as blue color in Fig. 5. The lower bounds of both $y_1 - y_2$ and $y_1 - y_3$ now become 0.2, making it infeasible to achieve $y_1 - y_2 \leq 0$ or $y_1 - y_3 \leq 0$. Therefore, we successfully do the verification with just one refinement process invocation.

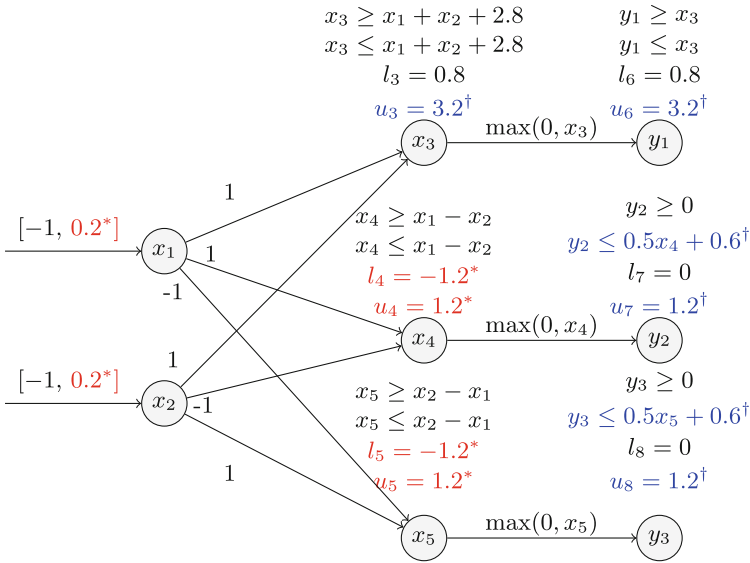


Fig. 5. The new intervals (in red with *) with two-adversarial labels encoding and new abstraction introduced by the second run of DeepPoly (in blue with †) (Color figure online)

3 Methodologies

As described in Sect. 2, we identify the feasible region of the network encoding and the negation of a property (i.e. $\Pi \wedge (y_c - y_t \leq 0) : t \neq c$) as a potential spurious region and leverage the refinement process to ascertain and possibly eliminate such spurious regions. To further improve the precision and efficiency, we propose three techniques as we have summarized in Sect. 1:

1. We update the negation of the property encoding to capture multiple spurious regions at the same time, as we demonstrate the example on $(y_1 - y_2 \leq 0) \vee (y_1 - y_3 \leq 0)$ in Sect. 2.2. This method allows us to reuse the linear programming part among several spurious regions and improve efficiency.
2. We leverage the multi-ReLU convex abstraction proposed in PRIMA [9] to obtain a more precise network encoding Π , which helps to increase the verification precision.
3. We detect adversarial examples to falsify robustness property. In particular, as the LP solver finds the conjunction of the network encoding and the negation of the property to be feasible, its optimization solution could actually ascertain a property violation and help us conclude with falsification.

We will discuss the three methodologies in separate subsections, and conclude this section with an overall description of our verification framework ARENA. Table 2 shows the notations we use in the main text of this section.

Table 2. Notations and descriptions of Sect. 3

| | |
|-----------|--|
| Π | The network constraint set/encoding |
| Ω | The multi-ReLU constraint set |
| Λ | The involved variable set during convex computation |
| Θ | The initial multidimensional octahedra during convex computation |
| P_i^H | The convex polytope in H-representation |
| P_i^V | The convex polytope in V-representation |

3.1 Multi-ReLU Network Encoding

As mentioned before, the constraint set subject to linear programming resolution is a conjunction of the network encoding and the negation of the property. In this subsection, we describe our network constraint construction. In the next subsection, we will describe the encoding of the negated property. Particularly, we capture the dependencies between the ReLU neurons in the same layer in our network encoding by leveraging the multi-ReLU convex relaxation in PRIMA.

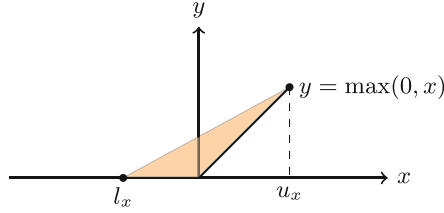


Fig. 6. The triangle approximation of a single ReLU neuron

As depicted in Fig. 6, DeepSRGR uses a triangular shape to encode each ReLU neuron independently, where the ReLU node $y = \max(0, x)$ with $x \in [l_x, u_x]$. The triangular shape is defined by three linear constraints:

$$y \geq x \qquad y \geq 0 \qquad y \leq \frac{u_x}{u_x - l_x}(x - l_x)$$

This looseness of ReLU encoding can inhibit precision improvement in DeepSRGR. As a matter of fact, it has been reported that, when they increase the maximum number of iterations from 5 to 20, only two more properties can be verified additionally, and no more properties can be verified when they further increase from 20 to 50 [10].

To break this precision barrier, we deploy the technique of *multi-ReLU relaxation* in PRIMA [9] where they compute the convex abstraction of k -ReLU neurons via novel convex hull approximation algorithms. For instance, if $k = 2$ and the ReLU neurons in the same layer are denoted by y_1, y_2 , and the inputs to these two ReLU neurons are x_1, x_2 respectively, PRIMA will compute a convex hull in (y_1, y_2, x_1, x_2) space to capture the relationship between the two ReLU neurons and their inputs. An example of the convex hull is defined as:

$$\Omega = \{ x_1 + x_2 - 2y_1 - 2y_2 \geq -2, \quad 0.375x_2 - y_2 \geq -0.75, \\ -x_1 + y_1 \geq 0, \quad -x_2 + y_2 \geq 0, \quad y_1 \geq 0, \quad y_2 \geq 0 \}$$

As we can see, Ω contains the constraint $x_1 + x_2 - 2y_1 - 2y_2 \geq -2$ that correlates (y_1, y_2, x_1, x_2) all together, which is beyond the single ReLU encoding. In general, PRIMA splits the input region into multiple sub-regions and then computes the convex hull of multiple ReLU neurons. For example, splitting the input region along $x_1 = 0$ results in two sub-regions where $y_1 = x_1$ (y_1 is activated) and $y_1 = 0$ (y_1 is deactivated). In each sub-region, the behavior of y_1 is determinate and this yields a tighter or even exact convex approximation. Finally, PRIMA computes a joint convex over-approximation (as in Ω) of the convex polytopes computed for each sub-region.

For deployment, we consider 3-ReLU neurons in our paper. We filter out the unstable ReLU neurons in each ReLU layer, and divide them into a set of 3-ReLU groups with one overlapping neuron between two adjacent groups as shown in Fig. 7, where a dashed box identifies a 3-ReLU group. We then leverage PRIMA to compute the constraints for each 3-ReLU group, and add those additional constraints into the original network encoding in order to obtain a more precise network abstraction and better verification precision.



Fig. 7. The 3-ReLU grouping for unstable ReLU neurons in the same layer i , where we use PRIMA to compute the convex relaxation for each group.

3.2 Multiple Adversarial Label Elimination

We now explain how we encode the negated property, especially when we take multiple spurious regions into consideration. As demonstrated in Sect. 2.2, to make it amenable for LP encoding, we need to compute the over-approximate convex hull of the union of multiple convex polytopes like in Fig. 4d. To explain the theory behind, we first introduce the required knowledge with respect to convex polytope representation. The convex polytope in this paper refers to a bounded convex polytope that is also a convex region contained in the n -dimensional Euclidean space R^n . There are two essential definitions of a convex polytope: as the intersection of half-space constraints (H-representation) and as the convex hull of a set of extremal vertices (V-representation) [14].

H-Representation. A convex polytope can be defined as the intersection of a finite number of closed half-spaces. A closed half-space in an n -dimensional space can be expressed by a linear inequality:

$$a_1x_1 + a_2x_2 + \dots + a_nx_n \leq b \quad (2)$$

A closed convex polytope can be taken as the set of solutions to a linear constraint set, just like the constraint set (1) shown in Sect. 2.2.

V-Representation. A closed convex polytope can also be defined as the convex hull with a finite number of points where this finite set must contain the set of extreme points of the polytope (i.e. the black-colored dots in Fig. 4d).

Double Description. The Double description method [12] aims to maintain both V-representation and H-representation during computation. This “duplication” is beneficial because to compute the intersection of two polytopes in H-representation is trivial since we only need to take the union set of the half-space constraints. On the other hand, to compute the convex hull of the union of two polytopes is trivial in V-representation as we take the union set of the vertices. The program `cddlib`² is an efficient implementation of the double description method, which provides functionalities that enable transformation from V-representation to H-representation (named as convex hull problem); and vice versa (named as vertex enumeration problem).

We leverage this V-H transformation in `cddlib` to compute the convex hull of the union of multiple convex polytopes. We set up a *batch size* δ to be in $[2, 5]$ ³, which defines the number of spurious regions to be considered simultaneously.

² <https://github.com/cddlib/cddlib>.

³ We explain our parameter range setting in Sect. 5 and also provide the batch size study experiments in Sect. 4.4.

Assume that the ground truth label is L_c , the related adversarial labels are L_1, \dots, L_δ , the convex-hull computation of $(y_c - y_1 \leq 0) \vee \dots \vee (y_c - y_\delta \leq 0)$ is conducted as follows:

Polytope Computation for Each Spurious Region. We compute the H-representation of the polytope for each spurious region in the $(\delta + 1)$ -dimensional space with respect to the variable set $\Lambda = y_c, y_1, \dots, y_\delta$. Intuitively, we obtain the H-representation of polytope $(y_c - y_i \leq 0)$ by taking the interval constraints of Λ (which is a multidimensional cube) conjunct with $y_c - y_i \leq 0$, as our example in Fig. 4. But this encoding is coarse as we neglect the dependencies between Λ that are in the same layer. For a more precise encoding, we follow the idea of [8] and compute the multidimensional octahedra Θ of $y_c, y_1, \dots, y_\delta$, which yields $3^{\delta+1} - 1$ constraints defined over Λ . Therefore, the H-representation of polytope $(y_c - y_i \leq 0)$ will be the constraint set Θ and $y_c - y_i \leq 0$.

Union of Convex Polytopes. We obtain the H-representation of the δ polytopes in the previous step and denote them by P_1^H, \dots, P_δ^H respectively. Since the union of polytopes is trivial in V-representation – as mentioned earlier, we use the H-V transformation in cddlib to generate these V-representations of the δ polytopes (referred to as P_1^V, \dots, P_δ^V). As illustrated in Fig. 8, we then produce the union set P_u^V of these vertices sets and transform it to its H-representation P_u^H , which is the convex hull of the union of δ polytopes. As P_u^H is represented by a set of linear inequalities, we conjunct it with the network encoding Π and submit the constraints for LP solving.

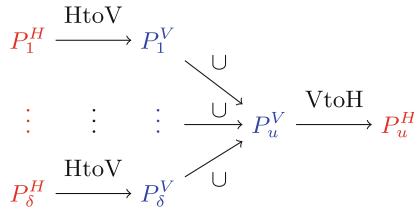


Fig. 8. The convex hull computation of the union of δ convex polytopes

3.3 Adversarial Example Detection

As mentioned previously, we take the conjunction of the network encoding and the negation of the property as the input constraint set to the LP solver and aim to eliminate spurious region(s) when detecting infeasibility. On the other hand, when the constraint set is feasible, we can set the input neurons and the unstable ReLU neurons as objective function and try to resolve for tighter intervals. In fact, a feasible constraint set indicates the possibility of a property violation. The LP solver not only returns the optimized value of the objective function, it also returns a solution that leads to the optimization, which could be a potential counter-example of robustness. Therefore, we include a supplementary procedure that takes each optimal solution obtained from the LP solver and checks if it constitutes an adversarial example.

This process brings forth two benefits: (1) it detects counter-examples and asserts the violation of robustness; (2) it enables the process to *terminate early with falsification* (once a counter-example is discovered) instead of exhausting all the iterations.

3.4 The Verification Framework ARENA

We now present an overview of our verification framework ARENA. In addition to the implementation of the three main technical points covered earlier, our framework includes the following optimizations as well.

Algorithm 1: Overall analysis procedure in ARENA

Input:

- N : input neural network with input layer γ_{in} , and output neurons y_1, \dots, y_n
- y_c : the output neuron corresponding to the ground truth label L_c ($1 \leq c \leq n$)
- δ is the refinement batch size (the number of adversarial labels to be eliminated by batch in each iteration).

Output: Verification result (Verified for robustness verified, Falsified for robustness violated, Inconclusive for inconclusive result).

```

1:  $(res, A_N) \leftarrow \text{VerifyByDeepPoly}(N)$  // Result and network abstraction
2: if  $res = \text{Verified}$  then
3:   return Verified
4: else
5:    $\Pi \leftarrow \text{GetConstraintsInNetwork}(N, A_N)$ 
6:    $\mathcal{L}_{adv} \leftarrow \{L_i \mid \text{IsFeasible}(\Pi \wedge y_c - y_i \leq 0), \forall i \neq c\}$  // All adversarial labels
7:    $\mathcal{L}_{adv} \leftarrow \text{SortByEstimatedSpuriousRegionSize}(\mathcal{L}_{adv})$  // Sort decreasingly
8:    $\mathcal{L}_{elim} \leftarrow \emptyset; i \leftarrow 0; iter\_num \leftarrow 999$  // Initialization
9:   while  $i < \text{Length}(\mathcal{L}_{adv})$  do
10:    if  $iter\_num > 2$  then
11:       $status, iter\_num \leftarrow \text{RefineWithKReLU}(N, \Pi, L_c, \mathcal{L}_{adv}[i], \mathcal{L}_{elim})$ 
12:      if  $status = \text{Falsified}$  or  $status = \text{Inconclusive}$  then
13:        return  $status$ 
14:      else //  $status = \text{Verified}$ 
15:         $\mathcal{L}_{elim} \leftarrow \mathcal{L}_{elim} \cup \{\mathcal{L}_{adv}[i]\}$ 
16:         $i \leftarrow i + 1$ 
17:    else
18:       $\mathcal{L}' \leftarrow \text{GetNextAdversarialLabelBatch}(\mathcal{L}_{adv}, \delta)$ 
19:       $status \leftarrow \text{EliminateAdversarialLabels}(N, \Pi, L_c, \mathcal{L}', \mathcal{L}_{elim})$ 
20:      if  $status = \text{Falsified}$  or  $status = \text{Inconclusive}$  then
21:        return  $status$ 
22:      else //  $status = \text{Verified}$ 
23:         $\mathcal{L}_{elim} \leftarrow \mathcal{L}_{elim} \cup \{\mathcal{L}'\}$ 
24:         $i \leftarrow i + \text{SizeOf}(\mathcal{L}')$ 
25:    if  $\mathcal{L}_{elim} = \mathcal{L}_{adv}$  then
26:      return Verified
27:    else
28:      return Inconclusive

```

Optimization 1: *Prioritising elimination of larger spurious regions.* We choose to order the sequence of the spurious regions according to the descending order of the respective regions’ sizes, by eliminating the “toughest” spurious region first. Since robustness only holds when all spurious regions are eliminated, we terminate the refinement process early if we fail to prune a larger spurious region. As it is difficult to compute the actual size of the spurious region, we deploy the metric in DeepSRGR where they take the lower bound of expression $y_c - y_i$ given by DeepPoly as the estimation of the region size, i.e. the smaller this value is, the larger the region is likely to be and thus it would be tougher for us to eliminate the region.

Optimization 2: *Cascading refinement.* Our system is designed to apply increasingly more scalable and less precise refinement methods. We hereby define process `RefineWithKReLU` as the refinement method with multi-ReLU encoding for the network and multi-adversarial label pruning feature *disabled*. Similarly, process `EliminateAdversarialLabels(δ)` is the refinement method with multi-ReLU encoding, and *taking into consideration* δ spurious regions *simultaneously*. With additional over-approximation error potentially being introduced by computing the union of polytopes, `EliminateAdversarialLabels(δ)` is less precise method compared to `RefineWithKReLU` but more scalable as it eliminates δ spurious regions simultaneously. We first use `RefineWithKReLU` to eliminate the larger spurious regions and record the number of iterations ς required to prune the current spurious region. If $\varsigma \leq 2$, this indicates that it is rather amiable to prune the current spurious region, and affordable to call upon `EliminateAdversarialLabels(δ)` to eliminate the remaining smaller spurious regions.

We present the overall analysis procedure in Algorithm 1. To begin with, we only apply the refinement process to images that fail to be verified by DeepPoly (lines 4). For refinement, we first obtain all the network constraints generated during abstract interpretation (line 5) and all potential adversarial labels (line 6). Then we call upon the processes `RefineWithKReLU` (line 11) or `EliminateAdversarialLabels(δ)` (line 19) to eliminate one or multiple spurious regions as stated in optimization 2 mechanism. The analyzer returns “Falsified” value if it detects an adversarial example (lines 13, 21); or it returns “Inconclusive” value if it fails to eliminate one of the adversarial labels and fails to find a counter-example (lines 13, 21, 28). We declare verification to be successful if and only if we can eliminate all the adversarial labels (lines 25–26).

Details of the two refinement processes are presented in Algorithms 2 and 3 (in Appendix B in our technical report). These two algorithms only differ in the property encoding (line 3–4 in Algorithm 3 vs lines 3–4 in Algorithm 2). Reading Algorithm 2 more closely: it first computes the convex hull of the union of the spurious regions according to Sect. 3.2 (line 3). Next, it conjuncts the convex hull with the network encoding in line 4, and add constraint $y_c - y_t > 0$ for each of the previously eliminated adversarial label L_t (lines 5–6); this helps to reduce the solution space further. If the combined constraint set Σ is found to be infeasible, the process returns “verified” since violation of the property cannot be attained (lines 7–8). But if Σ is feasible, the process leverages it to

further tighten the bounds for input and unstable ReLU neurons and updates the network (lines 9, 14). Moreover, the process checks if each LP solution is a valid counter-example; if so, it returns “Falsified” (lines 10–11, 15–16). With the newly solved bounds, the process then re-runs DeepPoly to obtain a tighter network encoding (lines 17–18) that is more amenable to encounter infeasibility in the latter iterations. Finally, the process returns “inconclusive” if it fails to conclude within the maximum number of iterations (line 20).

Algorithm 2: The refinement procedure `EliminateAdversarialLabels`

Function Name: `EliminateAdversarialLabels($N, \Pi, L_c, \mathcal{L}', \mathcal{L}_{elim}$)`
Input:

- N : input neural network with input layer γ_{in} , and output neurons y_1, \dots, y_n
- Π : the constraint set of N
- y_c : the output neuron corresponding to the ground truth label L_c ($1 \leq c \leq n$)
- $\mathcal{L}', \mathcal{L}_{elim}$: the batch of adversarial labels to be refined, and the list of previously eliminated labels

Output: the refinement status

```

1: counter = 0
2: while counter <  $\tau$  do //  $\tau$  is an iteration threshold
3:    $P_u^H \leftarrow \text{ComputeConvexHull}(N, \mathcal{L}')$  //  $P_u^H$  is the convex hull of polytopes
4:    $\Sigma \leftarrow \Pi \wedge P_u^H$  // Initialize constraint set
5:   for all  $L_t \in \mathcal{L}_{elim}$  do
6:      $\Sigma \leftarrow \Sigma \wedge (y_c - y_t > 0)$ 
7:   if IsInfeasible( $\Sigma$ ) then
8:     return Verified
9:    $N \leftarrow \text{LPSolveInputInterval}(\Sigma, \gamma_{in})$  // Update network with new bounds
10:  if ExistsAnAdversarialExample( $N$ ) then
11:    return Falsified
12:  for all ReLU layer  $\gamma_k$  in  $N$  do
13:     $\gamma_k \leftarrow \text{GetPrecedingInputAffineLayer}(\gamma_k')$ 
14:     $N \leftarrow \text{LPSolveUnstableReLUs}(\Sigma, \gamma_k)$  // Update new bounds
15:    if ExistsAnAdversarialExample( $N$ ) then
16:      return Falsified
17:   $A \leftarrow \text{RecomputeNetworkAbstractionByDeepPoly}(N)$ 
18:   $\Pi \leftarrow \text{GetConstraintsInNetwork}(N, A)$ 
19:  counter = counter + 1
20: return Inconclusive

```

4 Experiments

We implemented our method in a prototypical verifier called ARENA in both C and C++ programming languages (C++ is used for the k -ReLU computation feature, while the rest of the system is implemented in C). Our verifier is built on top of DeepPoly in [15]: it utilizes DeepPoly as the back-end abstract interpreter

for neural networks. Moreover, it uses **Gurobi**⁴ version 9.5 as the LP solver for constraints generated during abstract refinement.

We evaluate the performance of ARENA with state-of-the-art CPU-based incomplete verifiers including DeepSRGR [10], PRIMA [9], and DeepPoly [7]. Furthermore, we compare with a complete verifier α, β -CROWN [13], which is GPU-based and the winning tool of VNN-COMP 2022 [16]. The evaluation machine is equipped with two 2.40 GHz Intel(R) Xeon(R) Silver 4210R CPUs with 384 GB of main memory and a NVIDIA RTX A5000 GPU. The implementation is 64-bit based.

Note that DeepSRGR [10] was purely implemented in Python, while the main analysis in ARENA, PRIMA and DeepPoly were implemented in C/C++. Furthermore, this original version of DeepSRGR does not support convolutional networks nor the ONNX network format in our benchmark. Therefore, to avoid any runtime discrepancy introduced by different languages and to support our tested networks, we re-implemented the refinement technique of DeepSRGR in C, and conducted the experiment on the re-implemented DeepSRGR, where we release our re-implementation of DeepSRGR at this link: <https://github.com/arena-verifier/DeepSRGR>. The source code of our verifier ARENA is available online at: <https://github.com/arena-verifier/ARENA>.

4.1 Experiment Setup

Evaluation Datasets and Testing Networks. We chose the commonly used MNIST [17] and CIFAR10 [18] datasets. MNIST is an image dataset with handwritten digits, containing gray-scale images with 28×28 pixels. CIFAR10 includes RGB three-channel images with size 32×32 . Our testing image set consists of the first 100 images of the test set of each dataset, which is accessible from [15].

We selected fully-connected (abbreviated as FC) and convolutional (abbreviated as Conv) networks from [15] as displayed in Table 3, with up to around 50k neurons. We explicitly list the number of hidden neurons, the number of activation layers, trained defense⁵, and the number of candidate images for each network. Here, the candidate images refer to those testing images that can be correctly classified by the network and we only apply robustness verification on the candidate images.

Robustness Analysis. We conducted robustness analysis against L_∞ norm attack [19] with a perturbation parameter ϵ . Assuming that each pixel in the test image originally takes an intensity value p_i , it now takes an intensity interval $[p_i - \epsilon, p_i + \epsilon]$ after applying L_∞ norm attack with a specified constant ϵ .

This naturally forms an input space defined by $\times_{i=1}^n [p_i - \epsilon, p_i + \epsilon]$, and all the tools attempt to verify if all the “perturbed” images within the input space will be classified the same as the original image by the tested network. If so, we

⁴ <https://www.gurobi.com/>.

⁵ A trained defense refers to a defense method against adversarial samples, with the purpose of improving the robustness property of the network.

claim that the robustness property is verified. On the contrary, if we detect a counter-example with a different classification label, we assert the falsification of the robustness property. Finally, if we fail to conclude with verification or falsification, we return *unknown* to the user, meaning that the analysis result is inconclusive. We set up a challenging perturbation ϵ for each network and show in Table 3.

Table 3. Experimental fully connected and convolutional networks

| Network | Dataset | Type | ϵ | #Layer | #Neurons | Defense | Candidates |
|-------------|---------|------|------------|--------|----------|-------------|------------|
| M_3_100 | MNIST | FC | 0.028 | 3 | 210 | None | 98 |
| M_5_100 | MNIST | FC | 0.08 | 6 | 510 | DiffAI | 98 |
| M_6_100 | MNIST | FC | 0.025 | 6 | 510 | None | 99 |
| M_9_100 | MNIST | FC | 0.023 | 9 | 810 | None | 97 |
| M_6_200 | MNIST | FC | 0.016 | 6 | 1,010 | None | 99 |
| M_9_200 | MNIST | FC | 0.015 | 9 | 1,610 | None | 97 |
| M_convSmall | MNIST | Conv | 0.11 | 3 | 3,604 | None | 100 |
| M_convMed | MNIST | Conv | 0.1 | 3 | 5,704 | None | 100 |
| M_convBig | MNIST | Conv | 0.306 | 6 | 48,064 | DiffAI [20] | 95 |
| C_6_500 | CIFAR10 | FC | 0.0032 | 6 | 3,000 | None | 56 |
| C_convMed | CIFAR10 | Conv | 0.006 | 3 | 7,144 | None | 67 |

4.2 Comparison with the CPU-Based Verifiers

We present the robustness analysis results for ten networks in Table 3 and we describe the parameter configurations of our tool ARENA in Appendix C in our technical report. To execute PRIMA, we use the “refinepoly” domain in ERAN [15]. We report the experiment results drawn from different tools for MNIST and CIFAR10 networks in Table 4. For ARENA, we report the number of verified images, the number of falsified images and the average execution time for each testing image. DeepSRGR does not detect adversarial examples, neither does it attempt to assert violation of the property. PRIMA, on the other hand, returns two unsafe image for one MNIST network only (the detailed results and parameter setting are given in Appendix E in our report). Thus we omit the falsification column from the report for these two methods. Due to time limitation, for networks M_9_200, C_6_500 and C_convMed, we set a 2 h timeout for each image. If the refinement process fails to terminate before timeout, we consider the verification as inconclusive.

We observe from Table 4 that ARENA returns significantly more conclusive images (including both verified and falsified images) for all the networks than DeepSRGR, with comparable or even less execution time than that of DeepSRGR. ARENA also returns more conclusive images for all the networks than PRIMA, except for the subject MNIST_3_100, where ARENA returns less verified images than PRIMA. Our in-depth investigation reveals that it is because

Table 4. The number of verified/falsified images and average execution time (in seconds) per image for MNIST and CIFAR10 network experiments

| Neural Net | ARENA | | | DeepSRGR | | PRIMA | | DeepPoly | |
|-------------|--------|---------|--------|----------|--------|--------|--------|----------|-------|
| | Verify | Falsify | Time | Verify | Time | Verify | Time | Verify | Time |
| M_3_100 | 63 | 5 | 87.65 | 54 | 68.76 | 69 | 123.73 | 24 | 0.105 |
| M_5_100 | 77 | 7 | 250.39 | 67 | 153.75 | 53 | 19.15 | 25 | 0.522 |
| M_6_100 | 45 | 6 | 650.10 | 38 | 324.14 | 38 | 173.03 | 23 | 0.280 |
| M_9_100 | 44 | 10 | 1527.2 | 34 | 1004.4 | 34 | 191.60 | 30 | 0.587 |
| M_6_200 | 48 | 3 | 1514.2 | 35 | 1312.3 | 34 | 222.45 | 25 | 0.313 |
| M_9_200 | 43 | 6 | 3857.8 | 35 | 3536.7 | 29 | 238.63 | 29 | 0.536 |
| M_convSmall | 69 | 7 | 176.93 | 66 | 251.27 | 70 | 84.23 | 31 | 0.605 |
| M_convMed | 66 | 5 | 2054.9 | 60 | 2826.6 | 59 | 125.88 | 24 | 1.646 |
| C_6_500 | 31 | 9 | 2703.3 | 24 | 3985.2 | 20 | 269.96 | 16 | 12.22 |
| C_convMed | 31 | 7 | 3417.1 | 30 | 4385.4 | 30 | 230.74 | 21 | 3.87 |

two out of the three hidden layers have their ReLU neurons being encoded *exactly* with MILP in PRIMA.

These analysis results are better visualized in Fig.9 and 10 in Appendix D in [our technical report](#). As can be seen in Appendix D, ARENA generally returns more conclusive images than the rest of the tools. On average, ARENA returns 15.8% more conclusive images than DeepSRGR and 16.6% more conclusive images than PRIMA respectively for the testing networks. In summary, to the best of our knowledge, ARENA outperforms the current state-of-the-art approximated methods that run on CPU.

4.3 Comparison with the GPU-Based Verifier α, β -CROWN

Furthermore, we compare with the state-of-the-art tool α, β -CROWN (alpha-beta-CROWN) [13]. Note that this is a *complete* verification tool in the sense that it will produce a conclusive answer given sufficient amount of time.

We started our experiments with the version of α, β -CROWN available in August 2022. We report in detail here the average execution time for each image, the number of verified images and falsified images in Table 5 with five networks. We rerun our experiments with the availability of the November 2022 version of α, β -CROWN for all tested networks and present the results in Table 6. In terms of execution speed, we observe that α, β -CROWN is much superior to ARENA, mainly due to the deployment of GPU acceleration. In terms of the numbers of verified and falsified images, we note that ARENA can verify or falsify 79.3% of that of α, β -CROWN on average, for the upper seven subject tests. For the last four subject tests, α, β -CROWN introduces MIP encoding in their solution to capture *exact* ReLU functionality, and thus further enhancing the number of verified images. We are currently investigating techniques for implementing ARENA on GPU, with the goal to improve the number of verified/falsified images with reasonable time bound.

Table 5. The number of verified/falsified images and average execution time for ARENA and α, β -CROWN (version dated Aug 2022), time is presented in *seconds*

| Neural Net | ARENA | | | α, β -CROWN | | |
|-------------|--------|---------|--------------|------------------------|---------|--------------|
| | Verify | Falsify | Average time | Verify | Falsify | Average time |
| M_3_100 | 63 | 5 | 87.65 | 54 | 11 | 25.15 |
| M_5_100 | 77 | 7 | 250.3 | 53 | 10 | 48.67 |
| M_convSmall | 69 | 7 | 176.9 | 39 | 16 | 3.06 |
| M_convMed | 66 | 5 | 1625.9 | 30 | 18 | 2.90 |
| M_convBig | 53 | 30 | 589.52 | 49 | 24 | 4.21 |

Table 6. The number of verified/falsified images and average execution time of all tested networks for ARENA and α, β -CROWN (version dated Nov 2022), time is presented in *seconds*

| Neural Net | ARENA | | | α, β -CROWN | | |
|-------------|--------|---------|--------------|------------------------|---------|--------------|
| | Verify | Falsify | Average time | Verify | Falsify | Average time |
| M_3_100 | 63 | 5 | 87.65 | 81 | 13 | 40.81 |
| M_5_100 | 77 | 7 | 250.3 | 87 | 12 | 30.89 |
| C_6_500 | 31 | 9 | 2703.3 | 21 | 20 | 600.2 |
| C_convMed | 31 | 7 | 3417.1 | 36 | 22 | 283.7 |
| M_convSmall | 69 | 7 | 176.9 | 83 | 16 | 9.10 |
| M_convMed | 66 | 5 | 1625.9 | 82 | 18 | 7.24 |
| M_convBig | 53 | 30 | 589.52 | 60 | 29 | 109.17 |
| M_6_100 | 45 | 6 | 650.1 | 82 | 8 | 182.44 |
| M_6_200 | 48 | 3 | 1514.2 | 87 | 4 | 313.37 |
| M_9_100 | 44 | 10 | 1527.2 | 77 | 13 | 278.45 |
| M_9_200 | 43 | 6 | 3857.8 | 79 | 9 | 455.53 |

4.4 Multi-adversarial Label Parameter Study

In this experiment, we selected three networks from Table 3 to assess how the batch size parameter δ may impact the verification precision and execution time.

Theoretically, a larger value of δ may lead to a more efficient analysis process as it allows more adversarial regions to be eliminated at the same time. However, setting the parameter to be δ requires the computation of the union of δ convex polytopes. This in turn may introduce more over-approximation error and may jeopardize the analysis precision.

As δ aims to speed up the refinement process, we present the number of images that are verified through iterative refinement process and the average verification time for *those refined images*. As we only apply the refinement process to those testing images that DeepPoly fails to verify, the *refined images*

Table 7. The number of verified images through the refinement process (VTR) and average verification time per refined image for different δ setting.

| Network | ARENA | | | | | | | | | |
|---------|------------|---------|----------------|---------|----------------|---------|----------------|---------|----------------|---------|
| | (disabled) | | $(\delta = 2)$ | | $(\delta = 3)$ | | $(\delta = 4)$ | | $(\delta = 5)$ | |
| | VTR | Time(s) | VTR | Time(s) | VTR | Time(s) | VTR | Time(s) | VTR | Time(s) |
| M_3_100 | 41 | 142.91 | 39 | 144.67 | 39 | 135.62 | 39 | 124.99 | 38 | 127.77 |
| M_6_100 | 22 | 1414.4 | 22 | 1312.6 | 22 | 1250.5 | 22 | 1202.4 | 22 | 1047.9 |
| M_6_200 | 26 | 4809.7 | 23 | 2552.2 | 23 | 1828.2 | 23 | 1663.2 | 23 | 1297.0 |

refer to those images that are successfully verified through our refinement process, NOT through the original DeepPoly process. The experimental results are shown in Table 7 where we compare among parameters $\delta = 2, 3, 4, 5$ and with multi-adversarial label feature being disabled (the same as setting $\delta = 1$).

The experiment results show that the choice of a larger δ still allows us to achieve closely comparable precision while requires less execution time. Since an appropriate set-up of parameters leads to a better combination of precision and efficiency, we describe our configuration of each tested network in Appendix C in our technical report.

5 Discussion

We now discuss the limitation of our work. As described in Sect. 3.2, our batch size parameter δ is bounded to 5 at maximum for both precision and time-efficiency concern. In consideration for precision solely, as we compute the over-approximate convex hull of the union of multiple convex polytopes, the process will inevitably introduce additional over-approximate error into the LP encoding, yielding coarser neuron intervals. Thus we bound the value of δ to mitigate the degree of precision sacrifice. For time-efficiency issue, the transformation between V-representation and H-representation (refer to Sect. 3.2) – in either direction – is generally NP-hard, thus incurring exponential overhead with larger dimensions. As the parameter δ yields a $(\delta + 1)$ -dimensional space, it is advisable to keep δ -value small so that the convex hull computation process will not become an execution bottleneck. For future work, we will explore the possibility of assigning δ dynamically for different networks to strike a better trade-off between speed and precision.

Our proposed refinement process could be applied to other verification techniques for improved precision, as long as they use linear constraints to approximate the underlying network [6, 8, 9].

6 Related Work

Network verification methods can be generally categorized as complete or incomplete methods. Complete methods conduct exact analysis over the network,

especially for ReLU-activated networks. Given adequate time and resources, the complete methods return deterministic verification or violation of the robustness property to the user. Typical existing works are usually SMT (satisfiability modulo theory) based, MILP (mixed integer linear program) based or branch and bound (BaB) based [3, 4, 21, 22]. For instance, β -CROWN [22] is a GPU-based verifier which uses branch and bound method to enable exact reasoning over the ReLU activation function. Furthermore, β -CROWN could also perform as an incomplete verifier with early termination.

On the other hand, the incomplete methods choose to over-approximate the non-linearity of the network using abstract interpretation or bound propagation etc. [6, 7, 23, 24]. They are faced with precision loss due to the over-approximation of network behaviour. Consequently, the analysis result becomes inconclusive when the incomplete verifiers fail to verify the property. To rectify this deficiency, researchers have proposed various techniques like [8–10]. The work in [8] presents a new convex relaxation method that considers multiple ReLUs jointly in order to capture the correlation between ReLU neurons in the same layer. This idea has been further developed in PRIMA [9] which reduces the complexity of ReLU convex abstraction computation via a novel convex hull approximation algorithm. In comparison, DeepSRGR [10] elects to refine the abstraction in an iterative manner, where it repeatedly uses the spurious regions to stabilize the ReLU neurons until the abstraction is precise enough to eliminate the adversarial label linked to that specified spurious region. In our work, we combine both these refinement methods [9, 10] to break the precision barrier and also leverages the double-description method to retain efficiency as well.

7 Conclusion

We leverage the double description method in convex polytope area to compute the convex hull of the union of multiple polytopes, making it amenable for eliminating multiple adversarial labels simultaneously and boosting the analysis efficiency. Furthermore, we combine the convex relaxation technique with the iterative abstract refinement method to improve the precision in abstract interpretation based verification system. We implemented our prototypical analyzer ARENA to conduct both robustness verification and falsification. Experiment results show affirmatively that ARENA enhances abstract refinement techniques by attaining better verification precision compared to DeepSRGR, with reasonable execution time; it also competes favourably in comparison with PRIMA. Finally, it is also capable of detecting adversarial examples.

We believe that our proposed method can positively boost the effectiveness of sound but incomplete analyses and be applied to other methods that use linear constraints to approximate the network for effective precision enhancement.

Acknowledgement. This research is supported by a Singapore Ministry of Education Academic Research Fund Tier 1 T1-251RES2103 and the National Research Foundation, Singapore under its Emerging Areas Research Projects (EARP) Funding Initiative. Any opinions, findings and conclusions or recommendations expressed in this

material are those of the author(s) and do not reflect the views of National Research Foundation, Singapore. We are grateful to Julian R uth saraedum and Komei Fukuda for their prompt answer to our queries on cddlib. And we appreciate Mark Niklas M uller’s assistance to our queries on PRIMA.

References

1. Ren, K., Zheng, T., Qin, Z., Liu, X.: Adversarial attacks and defenses in deep learning. *Engineering* **6**(3), 346–360 (2020)
2. Yuan, X., He, P., Zhu, Q., Li, X.: Adversarial examples: attacks and defenses for deep learning. *IEEE Trans. Neural Netw. Learn. Syst.* **30**(9), 2805–2824 (2019)
3. Tjeng, V., Xiao, K., Tedrake, R.: Evaluating robustness of neural networks with mixed integer programming. In: International Conference on Learning Representations (ICLR). OpenReview.net (2019)
4. Katz, G., Barrett, C., Dill, D.L., Julian, K., Kochenderfer, M.J.: Reluplex: an efficient SMT solver for verifying deep neural networks. In: Majumdar, R., Kun ak, V. (eds.) CAV 2017. LNCS, vol. 10426, pp. 97–117. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-63387-9_5
5. Pulina, L., Tacchella, A.: An abstraction-refinement approach to verification of artificial neural networks. In: Touili, T., Cook, B., Jackson, P. (eds.) CAV 2010. LNCS, vol. 6174, pp. 243–257. Springer, Heidelberg (2010). https://doi.org/10.1007/978-3-642-14295-6_24
6. Gehr, T., Mirman, M., Drachler-Cohen, D., Tsankov, P., Chaudhuri, S., Vechev, M.: AI2: safety and robustness certification of neural networks with abstract interpretation. In: IEEE Symposium on Security and Privacy (SP), pp. 3–18. IEEE Computer Society (2018)
7. Singh, G., Gehr, T., P uschel, M., Vechev, M.T.: An abstract domain for certifying neural networks. *Proc. ACM Program. Lang.* **3**(POPL), 41:1–41:30 (2019)
8. Singh, G., Ganvir, R., P uschel, M., Vechev, M.T.: Beyond the single neuron convex barrier for neural network certification. In: Wallach, H.M., Larochelle, H., Beygelzimer, A., d’Alch e-Buc, F., Fox, E.B., Garnett, R., (eds.) *Advances in Neural Information Processing Systems*, vol. 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8–14, 2019, Vancouver, BC, Canada, pp. 15072–15083 (2019)
9. M uller, M.N., Makarchuk, G., Singh, G., P uschel, M., Vechev, M.T.: PRIMA: general and precise neural network certification via scalable convex hull approximations. *Proc. ACM Program. Lang.* **6**(POPL), 1–33 (2022)
10. Yang, P.: Improving neural network verification through spurious region guided refinement. In: Groote, J.F., Larsen, K.G. (eds.) TACAS 2021. LNCS, vol. 12651, pp. 389–408. Springer, Cham (2021). https://doi.org/10.1007/978-3-030-72016-2_21
11. Clarke, E., Grumberg, O., Jha, S., Lu, Y., Veith, H.: Counterexample-guided abstraction refinement. In: Emerson, E.A., Sistla, A.P. (eds.) CAV 2000. LNCS, vol. 1855, pp. 154–169. Springer, Heidelberg (2000). https://doi.org/10.1007/10722167_15
12. Fukuda, K., Prodon, A.: Double description method revisited. In: Deza, M., Euler, R., Manoussakis, I. (eds.) CCS 1995. LNCS, vol. 1120, pp. 91–111. Springer, Heidelberg (1996). https://doi.org/10.1007/3-540-61576-8_77

13. CMU. Alpha-Beta-CROWN: a fast and scalable neural network verifier with efficient bound propagation (2022). <https://github.com/huanzhang12/alpha-beta-CROWN>. Accessed 11 Aug 2022
14. McMullen, P.: Convex polytopes, by Branko Grunbaum, second edition (first edition (1967) written with the cooperation of V. L. Klee, M. Perles and G. C. Shephard. *Comb. Probab. Comput.* **14**(4), 623–626 (2005)
15. ETH. ETH Robustness Analyzer for Neural Networks (ERAN) (2022). <https://github.com/eth-sri/eran>. Accessed 11 Aug 2022
16. 3rd International Verification of Neural Networks Competition (VNN-COMP'22) (2022). <https://sites.google.com/view/vnn2022>. Accessed 11 Aug 2022
17. LeCun, Y., Cortes, C.: MNIST handwritten digit database (2010)
18. Krizhevsky, A., Nair, V., Hinton, G.: Cifar-10 (Canadian institute for advanced research)
19. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: IEEE Symposium on Security and Privacy (SP), pp. 39–57 (2017)
20. Mirman, M., Gehr, T., Vechev, M.T.: Differentiable abstract interpretation for provably robust neural networks. In: International Conference on Machine Learning (ICML), pp. 3575–3583 (2018)
21. Botoeva, E., Kouvaros, P., Kronqvist, J., Lomuscio, A., Misener, R.: Efficient verification of relu-based neural networks via dependency analysis. In: The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, 7–12 February 2020, pp. 3291–3299. AAAI Press (2020)
22. Wang, S.: Beta-crown: efficient bound propagation with per-neuron split constraints for complete and incomplete neural network verification. CoRR, abs/2103.06624 (2021)
23. Zhang, H., Weng, T.W., Chen, P.Y., Hsieh, C.J., Daniel, L.: Efficient neural network robustness certification with general activation functions. In: Bengio, S., Wallach, H.M., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R., (eds.) Advances in Neural Information Processing Systems Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3–8 December 2018, Montréal, Canada, vol. 31, pp. 4944–4953 (2018)
24. Wang, S., Pei, K., Whitehouse, J., Yang, J., Jana, S.: Formal security analysis of neural networks using symbolic intervals. In: Enck, W., Felt, A.P., (eds.) 27th USENIX Security Symposium, USENIX Security 2018, Baltimore, MD, USA, 15–17 August 2018, pp. 1599–1614. USENIX Association (2018)