Accelerating Adversarially Robust Model Selection for Deep Neural Networks via Racing

Matthias König¹, Holger H. Hoos^{1,2}, Jan. N. van Rijn¹

¹LIACS, Leiden University, The Netherlands ²Chair for AI Methodology, RWTH Aachen University, Germany h.m.t.konig@liacs.leidenuniv.nl, hh@aim.rwth-aachen.de, jan.n.van.rijn@liacs.leidenuniv.nl

Abstract

Recent research has introduced several approaches to formally verify the robustness of neural network models against perturbations in their inputs, such as the ones that occur in adversarial attacks. At the same time, this particular verification task is known to be computationally challenging. More specifically, assessing the robustness of a neural network against input perturbations can easily take several hours of compute time per input vector, even when using state-ofthe-art verification approaches. In light of this, it becomes challenging to select from a given set of neural network models the one that is best in terms of robust accuracy, *i.e.*, the fraction of instances for which the model is known to be robust against adversarial perturbations, especially when given limited computing resources.

To tackle this problem, we propose a racing method specifically adapted to the domain of robustness verification. This racing method utilises Δ -values, which can be seen as an efficiently computable proxy for the distance of a given input to a neural network model to the decision boundary. We present statistical evidence indicating significant differences in the empirical cumulative distribution of Δ -values for robust and non-robust instances. Using this information, we show that it is possible to reliably expose vulnerabilities in the model with relatively few input iterations. Overall, when applied to selecting the most robust network from sets of 31 MNIST and 27 CIFAR-10 networks, our proposed method achieves speedups of a factor of 108 and 42, respectively, in terms of cumulative running time compared to standard local robustness verification on the complete testing sets.

Introduction

Deep learning methods based on neural networks have gained increasing prominence in safety-critical domains and use contexts, such as unmanned mobile phone face recognition systems or aircraft manoeuvre advisory systems; see, *e.g.*, Julian, Kochenderfer, and Owen (2019). At the same time, neural networks are well known to be susceptible to adversarial attacks, where a slight modification of the input can lead to misclassification by the network (Szegedy et al. 2014). These perturbations can be so subtle that they remain imperceptible to human eyes, facilitating the need for the development of methods to formally reason about the behaviour of a neural network. In this context, formal verification methods have been developed that seek to assess the robustness of a trained neural network with respect to various input perturbations under specific norms, such as the commonly used l_{∞} -norm (Botoeva et al. 2020; Henriksen and Lomuscio 2020; Bunel et al. 2018; Dvijotham et al. 2018; Gehr et al. 2018; Wang et al. 2018a; Xiang, Tran, and Johnson 2018; Ehlers 2017; Katz et al. 2017; Tjeng, Xiao, and Tedrake 2019; Bastani et al. 2016; Papernot et al. 2016; Scheibler et al. 2015; Goodfellow, Shlens, and Szegedy 2015). Notice that formal verification methods can, in principle, verify a neural network model with respect to any pre-defined property describing the input-output behaviour of the model. In this work, we focus exclusively on verification with respect to adversarial input perturbations.

Generally, neural network verification is highly complex, and even simple network properties have been shown to be NP-complete problems (Katz et al. 2017). To solve these problems, state-of-the-art verification algorithms rely on the usage of sophisticated solvers, including mixed integer programming and satisfiability modulo theories solvers, and often demand several hours of running time to solve a single verification problem. This holds even for relatively small networks, such as those trained on the MNIST dataset (see, e.g., König, Hoos, and van Rijn 2022). Thus, much recent work has been concerned with the development of more efficient verification algorithms, e.g., by employing the branch-and-bound method for solving the verification problem (De Palma et al. 2021; Bunel et al. 2020; Wang et al. 2018a; Bunel et al. 2018; Ehlers 2017) or by tightening bounds in the problem formulation using symbolic interval propagation (Botoeva et al. 2020; Henriksen and Lomuscio 2020; Wang et al. 2018b, 2021) and abstraction (Bak et al. 2020; Singh et al. 2019b; Zhang et al. 2018; Gehr et al. 2018; Singh et al. 2018) techniques. However, even in light of recent developments, neural network verification remains a challenging and expensive computational task, especially as network complexity and dataset size increase.

Neural network verification can be divided into *local* and *global* verification (Sun et al. 2022). In this work, we focus on local robustness verification. Local robustness verification typically considers a trained neural network, along with a set of inputs and a verification property specification. Similar to other performance metrics of a neural net-

Copyright © 2024, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

work model, such as accuracy, one can compute *robust accuracy* by counting the fraction of inputs that are provably robust with regard to the given property. However, this adds significant overhead to the evaluation procedure, due to the high computational demands raised by most formal verification algorithms as explained earlier. Consequently, this overhead grows substantially if multiple models are considered and compared against each other in (robust) performance; a scenario not only faced by practitioners but also typically encountered during Neural Architecture Search (NAS), where the goal is to select a suitable model from a large search space (see, *e.g.*, Elsken, Metzen, and Hutter 2019). In this context, adding robust accuracy as a selection criterion would hardly be feasible due to the large computational costs.

In this work, we seek to improve the efficiency of local robustness verification from a previously unexplored, metaalgorithmic point of view. Specifically, we propose a method to efficiently evaluate and compare the robustness of different neural network models (or variations of the same model) against adversarial attacks. Moreover, we consider the problem of selecting the most robust model, *i.e.*, the model with the highest certified robust accuracy, from a given set of trained neural networks whilst making the most efficient use of the computational budget.

In summary, our proposed method employs a racing algorithm in which the considered neural network models are subjected to local robustness verification with respect to adversarial attacks. After each input iteration, their performance (in terms of robust accuracy) is measured and the verification procedure stops for a given network as soon as its robust accuracy is lower than the robust accuracy obtained by its competitors. Racing approaches are well studied and have already been successfully employed in other, resource intense domains, such as hyperparameter optimisation (Hutter, Hoos, and Leyton-Brown 2011; Birattari et al. 2010).

Complementary to the racing approach, we propose a novel sampling strategy based on the likelihood of a given input instance being adversarially robust. Essentially, this strategy prioritises input instances during the verification procedure that are most likely to expose vulnerabilities of the neural network model and, therefore, provide valuable insights into its robustness after fewer input iterations of the verification procedure and, hence, at a lower computational cost. At the same time, it reduces the risk of selecting suboptimal models, which might show higher robust accuracy than other candidate models after verifying some randomly sampled input instances but might perform worse overall. In fact, when using random sampling, the only way to mitigate this risk would be to increase the number of input iterations – with the associated costs involved.

To enable the proposed sampling strategy, we must define or estimate the likelihood of a network input being adversarially robust. In our case, this involves estimating their proximity to the decision boundaries of the model, captured by means of Δ -values, which will be explained in the following. Using this strategy, we can bias the sampling towards inputs for which adversarial attacks are most likely to occur. Although the relation between adversarial examples and the Algorithm 1: Racing approach for robust model selection

Input: Trained neural network models $\mathcal{N} = \{N_1, N_2, \dots, N_m\}$; Network input instances $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$; Verification algorithm VERIFY (N_i, x_j) that returns *sat*, *unsat* or *unsolved*;

Output: Model with highest robust accuracy: N_{selected}

1: $\mathcal{C} \leftarrow \mathcal{N}$ 2: $u_i \leftarrow 0$ with i = 1, 2, ..., |N|3: for all $x_i \in \mathcal{X}$ do for all $N_i \in \mathcal{C}$ do 4: 5: if VERIFY (N_i, x_j) is unsat then 6: $u_i \leftarrow u_i + 1$ end if 7: 8: end for $\mathcal{C} \leftarrow \{ N_i \mid i \in \arg\max_i \{ u_i \} \}$ 9: 10: end for 11: Randomly select one element N_{selected} from set C12: **Return** N_{selected}

decision boundary of a neural network has been extensively studied (Ding et al. 2020; Zhang et al. 2020; Croce, Andriushchenko, and Hein 2019; He, Li, and Song 2018; Liu et al. 2016), we are not aware of any existing work leveraging these insights in the context of local robustness verification procedures.

To our knowledge, the work presented here is the first to investigate techniques to increase the efficiency of adversarially robust model selection for neural networks. Specifically, it tackles the problem of selecting the most robust neural network model from a given set of models, whilst reducing the amount of compute time needed to obtain robustness certificates for the given input instances. To this end, we propose an efficient model selection method based on a novel heuristic that reliably quantifies the likelihood of a network being adversarially robust with respect to a given input. Concretely, we introduce Δ -values, which serve as a proxy for the distance from an input instance to the decision boundaries of a neural network model, and we present statistical evidence indicating significant differences in the empirical cumulative distribution of these Δ -values for robust and non-robust instances.

We considered two sets of models, each containing a large and diverse number of neural networks trained on the MNIST and CIFAR-10 datasets, respectively. For these sets, we demonstrate that our proposed model selection method reduces the cumulative running time required for selecting the most robust neural network model by a factor of 108 for the considered set of MNIST networks and a factor of 42 for the considered set of CIFAR networks when compared against (selection based on) exhaustive evaluation, where each model is verified with respect to all available input instances during the verification procedure.

Adversarially Robust Model Selection

The work presented here considers the selection of adversarially robust neural network models. More specifically, in the evaluation phase of several different neural network models, a model is to be determined that achieves the highest robust accuracy with respect to adversarial attacks.

Our proposed method has two main components: a racing approach and a sampling strategy based on a sorting mechanism for the input instances on which the network is verified. We considered two variants of the racing approach. The first one is a naïve racing approach in which the best-performing candidate models are selected at every input iteration, whereas the second one represents an adaptation of the *F-Race* algorithm (Birattari et al. 2010), which gathers statistical evidence against some candidate models before they are discarded. Both variants of the racing approach as well as our proposed sorting mechanism will be further explained in the following.

Naïve Racing Approach

Generally, the idea of a racing approach is to evaluate a finite set of candidate models while allocating the computational resources among them in a systematic way (see, *e.g.*, Maron and Moore 1993). To do so, the racing approach verifies step-by-step each candidate model in the given set, where in this context, a step corresponds to an input instance on which the neural network models are verified. At each step, all the remaining candidate models are verified, possibly in parallel, and candidate models are discarded once they are outperformed by others, *i.e.*, once one or more networks have obtained a higher robust accuracy.

An overview of this approach, which we refer to as the naïve racing approach for the remainder of this paper, can be found in Algorithm 1. After each iteration over the input instances, it identifies the model with the highest robust accuracy (determined based on u_i which represents the number of unsat instances for network N_i) and updates the set of candidate models C accordingly. Notice that the selection criterion on line 9 can by virtue of the arg max operator return a set of multiple networks. Moreover, u_i increases whenever a network N_i is found to be robust, *i.e.*, *unsat*, w.r.t. to a given input. On the other hand, an instance that is misclassified by the model would be considered as *sat*. The algorithm stops once all input instances have been processed, and the final output is the model with the highest determined robust accuracy.

F-Race

An important aspect of the model selection problem outlined above is that it can be viewed as a stochastic problem. In fact, although the process of formally verifying the behaviour of a neural network model with respect to certain input instances is deterministic (*i.e.*, multiple runs on the same input instance will always lead to the same result), its outcome depends on the particular instance to which it is applied. Concurrently, the specific instance being verified can be regarded as having been sampled from an underlying probability distribution, which may be unknown. For the naïve racing approach, this could lead to models being prematurely discarded after a few input iterations, even if that model would achieve the highest robust accuracy overall, *i.e.*, if it was verified with respect to all available input instances.

To address this stochasticity, Birattari et al. (2010) proposed F-Race, a widely known, state-of-the-art racing algorithm. F-Race can be considered an extension of the naïve racing approach, where the naïve selection criterion (line 9 in Algorithm 1) is replaced with a statistical test. Concretely, after each iteration over the input instances, F-Race performs a statistical test, typically, the non-parametric Friedman test, to determine if there are significant differences in the number of *unsat* instances per neural network model. If the null hypothesis is rejected or, in other words, significant differences exist, F-Race applies post-tests to identify the models which are performing statistically significantly worse than the best, and updates C accordingly. The algorithm stops when all input instances have been processed, and the final output is the model with the highest robust accuracy.

Since we are interested in the fraction of instances that are *unsat*, we used Cochran's Q test to determine if there are significant differences among the *unsat* counts for each of the networks. Notice that Cochran's Q is identical to the Friedman test but applicable when the responses are binary. When only two candidate networks remain, we used the McNemar test (without continuity correction), which can be seen as a special case of Cochran's Q test (Tate and Brown 1970). Any significant Cochran's Q (or McNemar) statistic is followed by Dunn's post-hoc test with a significance threshold of p = 0.05, and networks are selected if they have a significantly higher certified robust accuracy than their competitors.

Sorting Mechanism

In addition, we propose a sampling strategy based on a mechanism that sorts the considered input instances according to their likelihood of being adversarially robust. The key idea behind this mechanism is that by exposing a neural network model to inputs that are least likely to be adversarially robust, we can more quickly gather insights into its vulnerability or, similarly, its robustness. In other words, if we initially verify a neural network model on its most "challenging" input instances, *i.e.*, instances on which it is most likely not robust, but obtain robustness guarantees for these instances, we can at least heuristically assume the model to also be robust with respect to the remaining instances.

To enable this sorting mechanism, we must define or estimate the likelihood of a neural network model input being adversarially robust. In our case, this involves estimating their distance from the decision boundaries of the model, captured by means of network outputs. Intuitively, if an input lies very close to the boundary between two classes, it can be assumed that small perturbations, such as those applied to adversarial examples, have a higher chance to change the prediction made by the model.

In this work, we estimate the distance to an adjacent class boundary as the difference between the neural network output corresponding to the most likely class and that corresponding to the second-most likely class, and we refer to this difference as Δ . Formally, we define $\Delta := \max(\{y_1, y_2, \ldots, y_n\}) - \max(\{y_1, y_2, \ldots, y_n\}) \setminus \max(\{y_1, y_2, \ldots, y_n\}))$, where y_n refers to the network out-



Figure 1: Empirical cumulative probability distribution of normalised Δ -values for *sat*, *unsat* and unsolved instances for the considered MNIST and CIFAR networks, respectively. Notably, the plot shows a statistically significant difference between the empirical distribution functions of Δ -values for *sat* and *unsat* instances. Specifically, for both MNIST and CIFAR networks, *sat* instances generally have smaller Δ -values than *unsat* instances. Statistical significance is determined by means of a Kolmogorov–Smirnov test with a significance threshold of 0.05.

put for a given class n. Based on the resulting Δ -values, we can, for each neural network model individually, sort the input instances in an non-decreasing order, where the smaller the value of Δ , the closer we assume an instance to lie to an adjacent class boundary.

Experimental Setup

We compiled two sets of neural network models: one set consisting of 31 neural networks trained on the MNIST dataset and one set containing 27 neural networks trained on the CIFAR-10 dataset. All networks were taken from the repository of the ERAN verification system (Müller et al. 2021; Singh et al. 2019a,b; Singh and Gehr 2019; Singh et al. 2018) and greatly vary in terms of architecture, training method as well as robust accuracy. Details of the considered networks can be found in the supplementary material. We verified each network for local robustness with respect to the first 100 instances in the test set of the MNIST and CIFAR-10 datasets, respectively.

To verify the MNIST networks, we used the state-ofthe-art complete CPU-based verification algorithm VeriNet (Henriksen and Lomuscio 2020) with a perturbation radius of $\epsilon = 0.04$, which lies well within the range of commonly chosen values for ϵ when verifying networks trained on MNIST (Wu et al. 2022; Botoeva et al. 2020; Wang et al. 2018a). Verification queries ran with a time budget of 3600 seconds on a cluster of machines equipped with Intel Xeon E5-2683 CPUs with 32 cores, 40 MB cache size and 94 GB RAM, running CentOS Linux 7.

To verify the more challenging CIFAR networks, we used β -CROWN, a state-of-the-art complete GPU-accelerated

verification method (Wang et al. 2021). For these networks, we verified local robustness with $\epsilon = 0.008$, a value in line with commonly chosen values of ϵ for networks trained on CIFAR (see, *e.g.*, Müller et al. 2022). Again, all verification queries ran with a time budget of 3600 seconds on machines equipped with NVIDIA GeForce GTX 1080 Ti GPUs with 11 GB video memory. Overall, the verification of the CIFAR networks used in our study consumed 558 hours in GPU time, whereas the verification of the MNIST networks demanded a total of 1380 CPU hours.

We note that, although the verification algorithms presented above are complete, they were sometimes unable to solve an instance due to time or memory limitations; we report such instances as unsolved.

Empirical Results

In the following, we will compare our proposed selection method against the F-Race approach, the naïve racing approach as well as selection based on exhaustive evaluation. The latter represents the conceptually simplest baseline for selecting the most robust model from a given set of neural network models. Using this approach, each model is verified with respect to all available input instances during the verification procedure. At each input iteration, the candidate model with the highest certified robust accuracy is selected as the incumbent; *i.e.*, the model that would be returned if the process were terminated at the given iteration. Differently from the racing approaches, the exhaustive evaluation approach does not eliminate any candidate models during the selection process; therefore, when run to completion, it will always achieve a regret of zero.



Figure 2: Number of candidate models as determined by each method after every input iteration. For the three methods that do not use the sorting mechanism, the line represents the average number of candidate models over 200 random input orders, each with different random seeds, along with along with the respective 95% confidence intervals. Clearly, naïve racing, coupled with our proposed sorting mechanism, reduces the number of candidates after substantially fewer input iterations than other methods. Notice that selection based on exhaustive evaluation does not eliminate models from the set of candidates, which, therefore, does not decrease in size.

We evaluated each method in terms of cumulative running time and regret. The former describes the total running time consumed by the verification algorithm until all given input instances have been processed and the most robust model has been determined. Regret, in the context of model selection, describes the difference between the performance of the selected model and the performance of the best model that could have been chosen based on complete and perfect knowledge. In other words, it represents the loss incurred by selecting a sub-optimal model.

Formally, the regret R is defined as follows. Suppose we have a set of candidate models $C = \{C_1, C_2, \ldots, C_n\}$, and we want to select one model from this set based on certified robust accuracy. Let C_{best} be the best model in the set, *i.e.*, the model with the highest certified robust accuracy ra. Then, $R := ra(C_{\text{best}}) - ra(C_{\text{selected}})$, where $ra(C_{\text{best}})$ represents the robust accuracy of the best model and $ra(C_{\text{selected}})$ the robust accuracy of the selected model.

Local Robustness at the Decision Boundary

First of all, we investigated the relationship between the local robustness of a neural network model and the estimated distance of an input instance from the decision boundary of the model. More specifically, we examined the empirical cumulative probability distribution of Δ -values across all considered models, giving rise to 3100 individual verification problem instances for MNIST and 2800 for CIFAR. Remember that Δ -values serve as a proxy for the distance of an instance from the closest adjacent class boundary. We normalised these values per network under consideration. The empirical cumulative distribution of the Δ -values is visualised in Figure 1. Notice that some instances could not be verified due to timeouts or memory limitations; we show these instances as unsolved. The plots clearly show that *sat* instances, *i.e.*, instances for which an adversarial example could be found, tend to have a smaller Δ -value than those that are *unsat*, *i.e.*, robust. The difference in distributions is determined as statistically significant by means of a Kolmogorov–Smirnov test with a standard significance threshold of 0.05.

At the same time, Figure 1 shows that there exist instances, which are found to be sat despite having a relatively large Δ -value, *i.e.*, a Δ -value close to the end of the (normalised) range of values. Upon further investigation, we found that for MNIST, such instances occurred for 12 out of the 31 neural network models we considered and for 15 out of the 27 CIFAR networks. Notice that for these models, no instance was found to be robust, which indicates that large Δ -values can occur also for *sat* instances if a neural network model generally suffers from poor robustness. However, this observation does not affect the performance of our proposed selection method, as models which are non-robust with respect to any input instance would be discarded from the set of candidate models early in the selection process regardless of their Δ -value and, hence, the sorting of input instances. Notice that when removing these neural network models from the set, the difference in Δ -values between sat and unsat instances grows even larger; more details can be found in the supplementary material.



Figure 3: Regret achieved by the considered methods, where regret describes the difference between the performance of the selected model and the performance of the best model that could have been chosen given all available information. For methods not using the sorting mechanism, the regret was averaged over 200 random input orders, each with different random seeds, and is shown with a 95% confidence interval. The plots show that naïve racing, coupled with our proposed sorting mechanism, achieves optimal regret with fewer input iterations than other methods.

Evaluation of Our Proposed Selection Method

We evaluated our proposed selection method, naïve racing coupled with the sorting mechanism, in terms of cumulative running time and regret, and compared its performance against the following three baselines: (i) F-Race, (ii) naïve racing without a sorting mechanism and (iii) selection based on exhaustive evaluation. For methods that do not employ the sorting mechanism (*i.e.*, all baselines), we repeated the selection process 200 times, where each time the order of the input instances was based on a different random seed. We report the average running time over all runs, along with the respective 95% confidence intervals.

Figure 2a displays the size of the set of candidate networks trained on the MNIST dataset throughout the selection process. It shows that our proposed selection method reduces the number of candidate models after fewer iterations compared to each considered baseline. At the same time, for the exhaustive evaluation approach, the number of considered models remains constant, resulting in a larger number of queries that need to be performed at every input iteration.

As the number of candidate models reduces very quickly, it could be assumed that the aggressive nature of our selection method might lead to a sub-optimal outcome of the model selection process. We investigated this potential trade-off and show the results in Figure 3a. As can be seen, every method reached an optimal regret, indicating that the significant speed-up does not necessarily compromise on the quality of the selection process. However, we note that some of the MNIST networks were found to be fully robust. These are, consequently, always selected by any of the selection methods, even those that are more aggressive. Lastly, notice that F-Race eliminates candidate models based on statistical evidence, which can lead to models being selected that are less robust than others but where this difference is not found to be statistically significant at the given iteration.

We also tested our method on networks trained on the more challenging CIFAR dataset. Neural networks trained on this dataset are generally more difficult to verify than those trained on the MNIST dataset (Li et al. 2020). Figure 2b shows the size of the set of candidate CIFAR networks throughout the selection process. Again, we found that our proposed selection method eliminates candidate models after fewer iterations compared to other methods. Concurrently, the difference between the naïve racing approach with and without the sorting mechanism is much smaller than the difference observed on MNIST networks.

However, Figure 3b shows the advantage of the sorting mechanism: the naïve racing approach using the sorting mechanism very quickly converges towards an optimal regret, while other methods either require substantially more iterations or do not reach the optimum at all. In fact, on this set of models, the naïve racing approach without the sorting mechanism always resulted in a sub-optimal model choice. Overall, these results clearly demonstrate that our new method can effectively select the most robust model, and does so in a more efficient way than F-Race, which discards models only after it obtained statistical significance between the robust accuracy of the candidate models.

Lastly, we studied in more detail the efficiency of our method compared to the baselines we considered, in terms of regret achieved for a specific time budget. This is visualised in Figure 4a for MNIST networks and Figure 4b for



Figure 4: Regret as a function of cumulative running time for each of the considered methods. Running time represents wallclock time on the machine on which the experiments were carried out. For methods not using the sorting mechanism, the regret was averaged over 200 random input orders, each with different random seeds, and is shown with a 95% confidence interval. The plots show that naïve racing, coupled with our proposed sorting mechanism, achieves optimal regret while using substantially less compute time than other methods. Each line ends once a specific method has processed all given input instances.

CIFAR networks. Notably, these plots reveal that for both sets of models, our method selects the best-performing, *i.e.*, most robust model while demanding less compute time than any of the considered baselines, especially selection based on exhaustive evaluation. In fact, for MNIST networks, the cumulative running required to complete the selection process is reduced by several orders of magnitude, *i.e.*, a 108fold speedup factor, when compared to selecting based on exhaustive evaluation (1380.93 vs 12.83 hours). Furthermore, for CIFAR networks, our selection method achieved a 41-fold speedup compared to the exhaustive evaluation approach (558.44 vs 13.18 hours). Generally, this decrease in cumulative running time occurs because our selection method iteratively eliminates models from the set of candidates, subsequently reducing the number of verification queries in the following iterations, as previously explained. We note that the number of verification queries directly depends on the number of models, which decreases throughout the selection process.

These results highlight that our proposed selection method is well-suited for scenarios in which computing resources are limited, as it is likely to select, within any given amount of running time, models that are more robust than those determined by the baselines considered in our study.

Conclusions and Future Work

In this study, we have, for the first time, demonstrated the effectiveness of advanced model selection techniques in the context of neural network verification. Specifically, we studied the problem of selecting the most robust neural network model from a given set of models whilst reducing the compute time needed to obtain robustness certificates for the given input instances.

To enable our proposed selection method, we introduced a novel sorting mechanism based on the likelihood of an input instance being robust with respect to adversarial input perturbations. This likelihood is captured by means of Δ -values, and we present statistical evidence indicating significant differences in the empirical cumulative distribution of these values for robust and non-robust instances. Overall, our method advises on the allocation of computing resources required to perform local robustness verification towards adversarially robust models and can, in principle, be used in combination with any verification system.

We empirically evaluated our method on two diverse sets of 31 and 27 neural networks, trained on the MNIST and CIFAR-10 datasets, respectively. Our results clearly show that our proposed model selection method significantly reduces the cumulative running time required to select the most robust neural network model from these sets. Specifically, compared to the exhaustive evaluation approach, our method achieved a speedup factor of 108 for the set of MNIST networks and a speedup factor of 42 for the set of CIFAR networks while still selecting the most robust model.

In future work, we plan to apply our method to other verification tasks (*e.g.*, robustness verification under bias field perturbations), network architectures and datasets, and to perform a systematic analysis of the relationship between Δ -values and the robustness of neural network models. In addition, we are interested in the precise relationship between the Δ -value and the distance to the nearest decision boundary.

Acknowledgments

This research was partially supported by TAILOR, a project funded by EU Horizon 2020 research and innovation program under GA No. 952215. The authors would like to thank Hadar Shavit for the insightful discussions.

References

Bak, S.; Tran, H.-D.; Hobbs, K.; and Johnson, T. T. 2020. Improved Geometric Path Enumeration for Verifying ReLU Neural Networks. In *Proceedings of the 32nd International Conference on Computer Aided Verification (CAV 2020)*, 66–96.

Bastani, O.; Ioannou, Y.; Lampropoulos, L.; Vytiniotis, D.; Nori, A.; and Criminisi, A. 2016. Measuring Neural Net Robustness with Constraints. In *Advances in Neural Information Processing Systems 29 (NeurIPS 2016)*, 2613–2621. Birattari, M.; Yuan, Z.; Balaprakash, P.; and Stützle, T. 2010. F-Race and Iterated F-Race: An Overview. In Bartz-Beielstein, T.; Chiarandini, M.; Paquete, L.; and Preuss, M., eds., *Experimental Methods for the Analysis of Optimization Algorithms*, 311–336. Springer.

Botoeva, E.; Kouvaros, P.; Kronqvist, J.; Lomuscio, A.; and Misener, R. 2020. Efficient Verification of ReLU-based Neural Networks via Dependency Analysis. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence* (*AAAI-20*), 3291–3299.

Bunel, R.; Lu, J.; Turkaslan, I.; Torr, P. H. S.; Kohli, P.; and Kumar, M. P. 2020. Branch and Bound for Piecewise Linear Neural Network Verification. *Journal of Machine Learning Research*, 21: 42:1–42:39.

Bunel, R.; Turkaslan, I.; Torr, P.; Kohli, P.; and Mudigonda, P. K. 2018. A Unified View of Piecewise Linear Neural Network Verification. In *Advances in Neural Information Processing Systems 31 (NeurIPS 2018)*, 1–10.

Croce, F.; Andriushchenko, M.; and Hein, M. 2019. Provable Robustness of ReLU networks via Maximization of Linear Regions. In Chaudhuri, K.; and Sugiyama, M., eds., *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTATS 2019)*, volume 89 of *Proceedings of Machine Learning Research*, 2057–2066.

De Palma, A.; Bunel, R.; Desmaison, A.; Dvijotham, K.; Kohli, P.; Torr, P. H. S.; and Kumar, M. P. 2021. Improved Branch and Bound for Neural Network Verification via Lagrangian Decomposition. *arXiv preprint arXiv:2104.06718*. Ding, G. W.; Sharma, Y.; Lui, K. Y. C.; and Huang, R. 2020. MMA Training: Direct Input Space Margin Maximization through Adversarial Training. In *Proceedings of the 8th International Conference on Learning Representations (ICLR* 2020), 2057–2066.

Dvijotham, K.; Stanforth, R.; Gowal, S.; Mann, T. A.; and Kohli, P. 2018. A Dual Approach to Scalable Verification of Deep Networks. In *Proceedings of the 38th Conference on Uncertainty in Artificial Intelligence (UAI 2018)*, 550–559.

Ehlers, R. 2017. Formal Verification of Piece-Wise Linear Feed-Forward Neural Networks. In *Proceedings of the 15th International Symposium on Automated Technology for Verification and Analysis (ATVA 2017)*, 269–286. Elsken, T.; Metzen, J. H.; and Hutter, F. 2019. Neural architecture search: A survey. *The Journal of Machine Learning Research*, 20(1): 1997–2017.

Gehr, T.; Mirman, M.; Drachsler-Cohen, D.; Tsankov, P.; Chaudhuri, S.; and Vechev, M. 2018. AI2: Safety and Robustness Certification of Neural Networks with Abstract Interpretation. In *Proceedings of the 39th IEEE Symposium on Security and Privacy (IEEE S&P 2018)*, 3–18.

Goodfellow, I. J.; Shlens, J.; and Szegedy, C. 2015. Explaining and Harnessing Adversarial Examples. In *Proceedings* of the 3rd International Conference on Learning Representations (ICLR 2015), 1–11.

He, W.; Li, B.; and Song, D. 2018. Decision Boundary Analysis of Adversarial Examples. In *Proceedings of the 6th International Conference on Learning Representations (ICLR* 2018), 1–15.

Henriksen, P.; and Lomuscio, A. 2020. Efficient Neural Network Verification via Adaptive Refinement and Adversarial Search. In *Proceedings of the 24th European Conference on Artificial Intelligence (ECAI 2020)*, 2513–2520.

Hutter, F.; Hoos, H. H.; and Leyton-Brown, K. 2011. Sequential Model-Based Optimization for General Algorithm Configuration. In *Proceedings of the 5th International Conference on Learning and Intelligent Optimization (LION 5)*, 507–523.

Julian, K. D.; Kochenderfer, M. J.; and Owen, M. P. 2019. Deep Neural Network Compression for Aircraft Collision Avoidance Systems. *Journal of Guidance, Control, and Dynamics*, 42(3): 598–608.

Katz, G.; Barrett, C.; Dill, D. L.; Julian, K.; and Kochenderfer, M. J. 2017. Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks. In *Proceedings of the* 29th International Conference on Computer Aided Verification (CAV 2017), 97–117.

König, M.; Hoos, H. H.; and van Rijn, J. N. 2022. Speeding up neural network robustness verification via algorithm configuration and an optimised mixed integer linear programming solver portfolio. *Machine Learning*, 111(12): 4565– 4584.

Li, L.; Qi, X.; Xie, T.; and Li, B. 2020. Sok: Certified robustness for deep neural networks. *arXiv preprint arXiv:2009.04131*.

Liu, W.; Wen, Y.; Yu, Z.; and Yang, M. 2016. Large-Margin Softmax Loss for Convolutional Neural Networks. In Balcan, M.; and Weinberger, K. Q., eds., *Proceedings of the 33nd International Conference on Machine Learning (ICML* 2016), volume 48, 507–516.

Maron, O.; and Moore, A. 1993. Hoeffding races: Accelerating model selection search for classification and function approximation. *Advances in Neural Information Processing Systems 6 (NeurIPS 1993)*.

Müller, C.; Serre, F.; Singh, G.; Püschel, M.; and Vechev, M. 2021. Scaling Polyhedral Neural Network Verification on GPUs. In *Proceedings of Machine Learning and Systems 3 (MLSys 2021)*, 1–14.

Müller, M. N.; Brix, C.; Bak, S.; Liu, C.; and Johnson, T. T. 2022. The Third International Verification of Neural Networks Competition (VNN-COMP 2022): Summary and Results. *arXiv preprint arXiv:2212.10376*.

Papernot, N.; McDaniel, P.; Wu, X.; Jha, S.; and Swami, A. 2016. Distillation as a Defense to Adversarial Perturbations Against Deep Neural Networks. In *Proceedings of the 37th IEEE Symposium on Security and Privacy (IEEE S&P 2016)*, 582–597.

Scheibler, K.; Winterer, L.; Wimmer, R.; and Becker, B. 2015. Towards Verification of Artificial Neural Networks. In *Proceedings of the 18th Workshop on Methoden und Beschreibungssprachen zur Modellierung und Verifikation von Schaltungen und Systemen (MBMV 2015)*, 30–40.

Singh, G.; Ganvir, R.; Püschel, M.; and Vechev, M. 2019a. Beyond the Single Neuron Convex Barrier for Neural Network Certification. In *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, 1–12.

Singh, G.; and Gehr, T. 2019. Boosting Robustness Certification of Neural Networks. In *Proceedings of the 7th International Conference on Learning Representations (ICLR 2019)*, 1–12.

Singh, G.; Gehr, T.; Mirman, M.; Püschel, M.; and Vechev, M. 2018. Fast and Effective Robustness Certification. In *Advances in Neural Information Processing Systems 31* (*NeurIPS 2018*), 1–12.

Singh, G.; Gehr, T.; Püschel, M.; and Vechev, M. 2019b. An Abstract Domain for Certifying Neural Networks. In *Proceedings of the 46th ACM SIGPLAN Symposium on Principles of Programming Languages (ACMPOPL 2019)*, 1–30.

Sun, W.; Lu, Y.; Zhang, X.; and Sun, M. 2022. DeepGlobal: A framework for global robustness verification of feedforward neural networks. *Journal of Systems Architecture*, 128: 102582.

Szegedy, C.; Zaremba, W.; Sutskever, I.; Bruna, J.; Erhan, D.; Goodfellow, I. J.; and Fergus, R. 2014. Intriguing properties of neural networks. In *Proceedings of the 2nd International Conference on Learning Representations (ICLR 2014)*, 1–10.

Tate, M. W.; and Brown, S. M. 1970. Note on the Cochran Q test. *Journal of the American Statistical Association*, 65(329): 155–160.

Tjeng, V.; Xiao, K.; and Tedrake, R. 2019. Evaluating Robustness of Neural Networks with Mixed Integer Programming. In *Proceedings of the 7th International Conference on Learning Representations (ICLR 2019)*, 1–21.

Wang, S.; Pei, K.; Whitehouse, J.; Yang, J.; and Jana, S. 2018a. Efficient Formal Safety Analysis of Neural Networks. In Advances in Neural Information Processing Systems 31 (NeurIPS 2018), 6369–6379.

Wang, S.; Pei, K.; Whitehouse, J.; Yang, J.; and Jana, S. 2018b. Formal Security Analysis of Neural Networks Using Symbolic Intervals. In *Proceedings of the 27th USENIX Security Symposium (USENIX Security 18)*, 1599–1614.

Wang, S.; Zhang, H.; Xu, K.; Lin, X.; Jana, S.; Hsieh, C.-J.; and Kolter, Z. 2021. Beta-CROWN: Efficient Bound Propagation with Per-neuron Split Constraints for Neural Network Robustness Verification. In Advances in Neural Information Processing Systems 34 (NeurIPS 2021), 29909–29921.

Wu, H.; Zeljic, A.; Katz, G.; and Barrett, C. W. 2022. Efficient Neural Network Analysis with Sum-of-Infeasibilities. In Fisman, D.; and Rosu, G., eds., *Proceedings of the 28th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS 2022)*, volume 13243, 143–163.

Xiang, W.; Tran, H.-D.; and Johnson, T. T. 2018. Output Reachable Set Estimation and Verification for Multilayer Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems*, 29(11): 5777–5783.

Zhang, H.; Weng, T.-W.; Chen, P.-Y.; Hsieh, C.-J.; and Daniel, L. 2018. Efficient Neural Network Robustness Certification with General Activation Functions. *Advances in Neural Information Processing Systems 31 (NeurIPS 2018)*, 31: 4944–4953.

Zhang, J.; Xu, X.; Han, B.; Niu, G.; Cui, L.; Sugiyama, M.; and Kankanhalli, M. S. 2020. Attacks Which Do Not Kill Training Make Adversarial Learning Stronger. In *Proceedings of the 37th International Conference on Machine Learning (ICML 2020)*, volume 119 of *Proceedings of Machine Learning Research*, 11278–11287.