# The Second International Verification of Neural Networks Competition (VNN-COMP 2021): Summary and Results

Stanley Bak[*], Changliu Liu[†], Taylor Johnson[‡]

**Abstract**

This report summarizes the second International Verification of Neural Networks Competition (VNN-COMP 2021), held as a part of the 4th Workshop on Formal Methods for ML-Enabled Autonomous Systems that was collocated with the 33rd International Conference on Computer-Aided Verification (CAV). Twelve teams participated in this competition. The goal of the competition is to provide an objective comparison of the state-of-the-art methods in neural network verification, in terms of scalability and speed. Along this line, we used standard formats (ONNX for neural networks and VNNLIB for specifications), standard hardware (all tools are run by the organizers on AWS), and tool parameters provided by the tool authors. This report summarizes the rules, benchmarks, participating tools, results, and lessons learned from this competition.

## 1   Introduction

Methods based on machine learning are increasingly being deployed for a wide range of problems, including recommendation systems, machine vision and autonomous driving. While machine learning has made significant contributions to such applications, few tools provide formal guarantees about the behaviours of neural networks.

In particular, for data-driven methods to be usable in safety-critical applications, including autonomous systems, robotics, cybersecurity, and cyber-physical systems, it is essential that the behaviours generated by neural networks are well-understood and can be predicted at design time. In the case of systems that are learning at run-time it is desirable that any change to the underlying system respects a given safety-envelope for the system.

While the literature on verification of traditionally designed systems is wide and successful, there has been a lack of results and efforts in this area until recently. The International Verification of Neural Networks Competition (VNN-COMP) was established in 2020, aiming to bring together researchers working on techniques for the verification of neural networks. In 2021, VNN-COMP[1] was held as a part of the 4th Workshop on Formal Methods for ML-Enabled Autonomous Systems (FoMLAS) that was collocated with the 33rd International Conference on Computer-Aided Verification (CAV).

While the first VNN-COMP in 2020 was a friendly competition where the participants tested their tools and reported the results in parallel, the second VNN-COMP in 2021 aims to provide a fair comparison and standardize the pipeline of the competition. Such standardization includes 1) standard formats where we use ONNX for neural networks and VNNLIB for specifications; and 2) standard hardware where all tools are run by the organizers on AWS either on CPU instances or cost-equivalent GPU instances. The competition was kicked off in January and the solicitation for participation was sent in February 2021. By March, several teams registered

---

[*]S. Bak is with Stony Brooks, `stanley.bak@stonybrook.edu`.

[†]C. Liu is with Carnegie Mellon University, `cliu6@andrew.cmu.edu`.

[‡]T. Johnson is with Vanderbilt University, `taylor.johnson@vanderbilt.edu`.

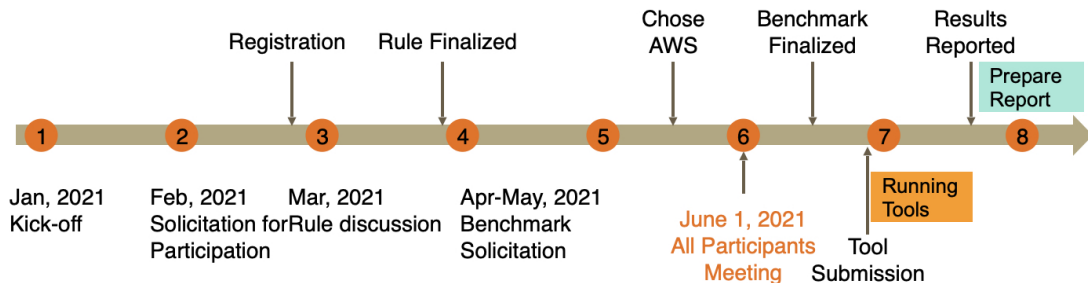[1]https://sites.google.com/view/vnn2021/home

Figure 1: Timeline for VNN-COMP2021.

to participate in this competition. The rule discussion was finalized in March 2021 where the finalized rules are summarized in section 2. From April to May 2021, the benchmarks to test the tools were solicited. Meanwhile, after comparing different choices, the organizing team finally decided to use AWS as our testing platform. On June 1st, 2021, all participants had an online meeting to agree on the rules and benchmarks. By the end of June 2021, twelve teams submitted their tools and the organizers spent two weeks running the tools on AWS to obtain the final results. The final results were reported in FoMLAS on July 19, 2021. The timeline is summarized in fig. 1. Most discussions took place on the repository https://github.com/stanleybak/vnncomp2021. There are three issues: rules discussion, benchmarks discussion, and tool submission, serving as the venues for the corresponding discussions. Moreover, the repository hosts all the submitted benchmarks and the scripts to run the tests.

The remainder of this report is organized as follows. Section 2 discusses the competition rules. Section 3 lists all participating tools, Section 4 lists all benchmarks, and Section 5 summarizes the results. Section 6 concludes the report and discusses potential future improvements.

## 2  Rules

**Terminology**  An *instance* is defined as (specification (pre- and post-condition), network, timeout). For example: an MNIST classifier with one input image, a given local robustness threshold $\epsilon$, and a specific timeout. A *benchmark* is defined as a set of instances. For example: a specific MNIST classifier with 100 input images, a given robustness threshold $\epsilon$, and one timeout per input.

**Run-time caps**  Per instance: any verification instance will timeout after at most X minutes, determined by the benchmark proposer. These can be different for each instance. Per benchmark: there is an upper runtime limit of 6 hours per benchmark. For example, a benchmark proposal could have six instances with a one hour timeout, or 100 instances each with a 3.6 minute timeout. To provide a fair comparison, we quantify the startup overhead for each tool by running it on small networks; and then we subtract the overhead from the total runtime.

**Instance score**  Each instance is scored is as follows:

- Correct hold: 10 points;

- Correct violated (where random tests or simple adversarial example generation did not succeed): 10 points;

- Correct violated (where random tests or simple adversarial example generation succeeded): 1 point;

- Incorrect result: -100 points.

**Time bonus** Time bonus is computed as follows.

- The fastest tool for each solved instance will receive +2 points.

- The second fastest tool will receive +1 point.

All runtimes below 1.0 seconds after overhead correction (explained below) are considered to be 1.0 seconds exactly for scoring purposes. If two tools have runtimes within 0.2 seconds (after all corrections), for scoring purposes we will consider them the same runtime.

**Overhead Correction** According to the rules discussion, we decided to subtract tool overhead time from the results. For example, simply importing tensorflow from Python and acquiring the GPU can sometimes take about 5 seconds, which would be unfortunate for benchmarks like ACASXu where some verification times are under a second.

To subtract overhead, we created trivial network instances and included those in the measurements. We then observed the minimum verification time along all instances, and considered that to be the overhead time for the tool.

One issue with that was that some tools had different overhead depending on if they were run in CPU mode or GPU mode, and this type of measurement penalized the GPU mode unintentionally. In the score reporting, we include a multi-overhead result where we apply the overhead measured for the mode the tool was actually run in.

**Benchmark score** The benchmark score of each category is a percentage. It is computed as 100 times the sum of the individual instance scores for that benchmark category divided by the maximum sum of instance scores of any tool for that benchmark category. For example, the tool with the highest sum of instance scores for a category should get 100%.

**Format** This year we standardized the inputs to be `onnx` neural networks and `vnnlib` specification files. Tool authors were also required to provide scripts to install their tool as much as possible, as well as run their tool on a specific instance provided the network file, specification file, and timeout. Specifications included simple disjunctions in both pre- and post-conditions to encode properties like an unchanged classification for an input in one of multiple hyperboxes. The specification is regarded as encoding a counter-example, meaning that a property is proven "correct" if the specification is shown to be unsatisfiable, while the property is shown to be violated, if a counterexample fulfilling the specification is found. Hence, robustness with respect to inputs in a hyper-box has to be encoded as disjunctive property, where any of the other classes is constrained to become the maximum output.

# 3 Participants

The following tools and teams participated in VNN-COMP. They are summarized in table 1.

## 3.1 Marabou

**Team** Guy Amir[1], Clark Barrett[2], Ahmed Irfan[3], Guy Katz[1], Teruhiro Tagomori[4], Alex Usvyatsov[1], Haoze Wu[2], Alexsandar Zeljic[2].
    [1] Hebrew University of Jerusalem, [2] Stanford University, [3] Amazon Web Services, [4] NRI Secure.

| Tool | GPU? | Floating Point Accuracy | Use of External Solvers |
|---|---|---|---|
| Marabou | - | LP-default | Gurobi |
| VeriNet | - | Either | Xpress Solver |
| ERAN | Yes | Sound | Gurobi |
| $\alpha,\beta$-CROWN | Yes | Either | Gurobi |
| DNNF | - | Either | None |
| NNV | - | 64bit | Matlab |
| OVAL | Yes | Either | None |
| RPM | - | 64bit | None |
| NV.jl | - | 64bit | None |
| Venus | - | 64bit | Gurobi |
| Debona | - | 32bit | Gurobi |
| nnenum | - | Either | None |

Table 1: Summary of key features of participating tools. For the "GPU" column, we only mark the tools that were tested on GPU during this competition. Tools that support GPU but were tested on CPU during this competition are not marked. For the floating point accuracy, "either" means the tool can handle either 32 or 64 bit floating point, depending on the onnx network; "LP-default" means the tool uses the default settings for LP solver; "sound" means the tool is fully floating-point-sound arithmetic with respect to IEEE-754 semantics up to the LP solver (up to 64x slower than non-fp-sound 32 bit arithmetic).

**Description** Marabou [22] is a Neural Network Verification toolkit that can answer queries about a network's properties by encoding and solving these queries as constraint satisfaction problems. It can accommodate networks with different activation functions (including ReLU, Leaky ReLU, Sign [1], Max, Absolute value) and topologies (e.g., FFNNs, CNNs, residual connections). It also uses the Split-and-Conquer algorithm [51] for parallelization to further enhance scalability. Marabou accepts multiple input formats, including protocol buffer files generated by the popular TensorFlow framework for neural networks, and the ONNX format.

The core of Marabou resolves around the Reluplex procedure [21], but it also supports multiple new techniques and solving modes. In particular, it incorporated the DeepPoly analysis introduced in [38] and the (MI)LP-based bound tightening first seen in [42]. In terms of complete verification procedure, in addition to the Reluplex procedure, Marabou also supports solving the verification query with a MILP encoding.

For the competition, Marabou uses the DeepPoly analysis to tighten the variable bounds and then uses a portfolio strategy for complete verification, with a fraction of the CPUs solving a MILP-encoding of the verification problem with Gurobi, and the rest running a new complete verification procedure that is currently under submission.

**Link** https://github.com/NeuralNetworkVerification/Marabou

**Commit** For reproduction of VNN-Comp results, use https://github.com/anwu1219/Marabou_private/commit/81e9f14f7ae9f6a2097524ea1291e86434ef42dc.

**Hardware and licenses** CPU, Gurobi License.

**Participated benchmarks** ACASXu, cifar10_resnet, eran, marabou-cifar10, mnistfc, oval21, verivital.

## 3.2 VeriNet

**Team** Patrick Henriksen, Alessio Lomuscio (Imperial College London).

**Description** VeriNet [17, 18] is a complete Symbolic Interval Propagation (SIP) based verification toolkit for feed-forward neural networks. The underlying algorithm utilises SIP to create a linear abstraction of the network, which, in turn, is used in an LP-encoding to analyse the verification problem. A branch and bound-based refinement phase is used to achieve completeness.

VeriNet implements various optimisations, including a gradient-based local search for counterexamples, optimal relaxations for Sigmoids, adaptive node splitting [17], succinct LP-encodings, and a novel splitting heuristic that takes into account indirect effects splits have on succeeding relaxations [18].

VeriNet supports a wide range of layers and activation functions, including Relu, Sigmoid, Tanh, fully connected, convolutional, max and average pooling, batch normalisation, reshape, crop and transpose operations, as well as additive residual connections.

Note that VeriNet subsumes the Deepsplit method presented in [18].

**Link** Scheduled for release late August 2021 at: https://vas.doc.ic.ac.uk/software/neural/.

**Hardware and licences** CPU and GPU, Xpress Solver license required for large networks.

**Participated benchmarks** ACASXu, `cifar10_resnet`, `cifar2020`, `eran`, `marabou-cifar10`, `mnistfc`, `nn4sys`, `oval21`, `verivital`.

## 3.3 ERAN

**Team** Mark Niklas Müller (ETH Zurich), Gagandeep Singh (UIUC), Markus Püschel (ETH Zurich), Martin Vechev (ETH Zurich)

**Description** ERAN [30, 33, 36–39] is a neural network verifier based on leveraging abstract interpretations to encode a network, pre- and post-condition as an LP or MILP problem. ERAN supports both incomplete and complete verification and can handle fully-connected, convolutional, and residual network architectures containing ReLU, Sigmoid, Tanh, and Maxpool non-linearities. It uses 64-bit precision (up to 16 times slower than 32-bit on some GPUs) and an arithmetic which is floating-point-sound (performs 4 times more operations) with respect to IEEE-754 semantics up to the LP solver. Single- and multi-neuron relaxations of non-linear activations [30] computed using GPUPoly [33] are combined with a partial MILP encoding and neuron-wise bound-refinement [39] to obtain a precise network encoding. The analyzer is written in Python, uses ELINA [40] for numerical abstractions, and Gurobi for solving LP and MILP instances. While the use of both GPU and CPU enables ERAN to utilize all resources of a system it also requires both a relatively strong GPU and CPU to be available in order to avoid one bottlenecking the other (which occurred with the GPU AWS instance).

When run in complete mode, ERAN generates concrete counterexamples. In incomplete mode, ERAN attempts to falsify a property by running a PGD attack before attempting verification using increasingly more expensive and precise abstractions. As we use the same abstractions for all instances of a benchmark, ERAN might fail to verify a property before exceeding the timeout. In these cases, a more expensive abstraction might have been able to verify the property in time.

**Link** https://github.com/eth-sri/eran

**Commit** Please use the main repository for anything but reproducing VNN-COMP results. https://github.com/mnmueller/eran_vnncomp2021.git
(1e474c77f72f86f450df9f0a860b4d35c490ea7c)

**Hardware and licences** CPU and GPU, GUROBI License

**Participated benchmarks** ACASXu, cifar10_resnet, cifar2020, eran, marabou-cifar10, mnistfc, nn4sys, oval21, verivital.

## 3.4 $\alpha,\beta$-CROWN

**Team** Huan Zhang[*] (Carnegie Mellon), Kaidi Xu[*] (Northeastern), Shiqi Wang[*] (Columbia), Zhouxing Shi (UCLA), Yihan Wang (UCLA), Xue Lin (Northeastern), Suman Jana (Columbia), Cho-Jui Hsieh (UCLA), Zico Kolter (Carnegie Mellon); * indicates equal contribution.

**Description** The $\alpha,\beta$-CROWN (alpha-beta-CROWN) verifier is based on an efficient bound propagation algorithm, CROWN [55], with a few crucial extensions [49, 53, 54]. We use the generalized version of CROWN in the auto_LiRPA library [53] which supports general neural network architectures (including convolutional layers, residual connections, recurrent neural networks and Tranformers) and a wide range of activation functions (e.g., ReLU, tanh, sigmoid, max pooling and average pooling), and is efficiently implemented on GPUs. We jointly optimize intermediate layer bounds and final layer bounds using gradient ascent (referred to as $\alpha$-CROWN or optimized CROWN/LiRPA [54]). Additionally, we use branch and bound [10] (BaB) and incorporate split constraints in BaB into the bound propagation procedure efficiently via the $\beta$-CROWN algorithm [49]. The combination of efficient, optimizable and GPU accelerated bound propagation with BaB produces a powerful and scalable neural network verifier.

Our verifier also utilizes a mixed integer programming (MIP) solver (Gurobi) for networks where MIP runs relatively fast, following the formulation in [42]. We use MIP to solve the tightest possible intermediate layer bounds for as many neurons as possible on CPUs within the timeout budget, and use $\alpha$-CROWN to solve the remaining ones on GPUs. Finally, we conduct BaB with $\beta$-CROWN using tightened bounds. Although the GPU AWS instance has weak CPUs, we still find that MIP is helpful for some benchmarks, and it can become more beneficial on a machine with both strong CPUs and GPUs. Note that $\alpha$-CROWN can exceed the power of a typical LP verifier when intermediate layer bounds are jointly optimized [32,54], so we do not use a LP solver to tighten bounds.

**Link** https://github.com/huanzhang12/alpha-beta-CROWN

**Commit** c12e6eeaf6b16f99a99b65f377d0f450d6466a83 (only for reproducing competition results; please use the main branch version for other proposes)

**Hardware and licenses** CPU and GPU with 32-bit or 64-bit floating point; Gurobi license required for mnistfc, eran, marabou-cifar10 and verivital benchmarks.

**Participated benchmarks** ACASXu, cifar10_resnet, cifar2020, eran, marabou-cifar10, mnistfc, nn4sys, oval21, verivital.

## 3.5 DNNF

**Team** David Shriver, Sebastian Elbaum, Matt Dwyer (University of Virginia).

**Description** DNNF [35] is a tool for neural network property falsification. It only attempts to find counter-examples to a property specification, and will not prove that a property holds. DNNF reduces properties and networks to robustness problems which can then be falsified using many different adversarial attack methods. Our reduction approach enables us to support much more complex property and network specifications, such as specifications with direct input and output relations, such as the relation $y > x$. The backend falsification method used for VNN-COMP is a custom implementation of PGD, but DNNF can also be run with several off-the-shelf

methods from foolbox or cleverhans, two popular python packages for adversarial attacks.

DNNF makes use of the DNNV framework [34] to load networks and properties, as well as to perform network simplifications. The current implementation of DNNF supports many different network operations. In particular it supports ONNX models with the following operations: Add, Atan, AveragePool, BatchNormalization, Concat, Conv, ConvTranspose, Elu, Flatten, Gather, Gemm, LeakyRelu, MatMul, MaxPool, Mul, Relu, Reshape, Resize, Shape, Sigmoid, Softmax, Sub, Tanh, Transpose, Unsqueeze.

DNNF can also be run on GPUs, which can speed up falsification for large models.

**Link** https://github.com/dlshriver/DNNF

**Commit** VNN-COMP results can be reproduced with commit `d4f08b43e4ad622157c65ac071183a3a0f4e6fe0`. For other uses, we suggest the `main` branch.

**Hardware and licences** DNNF can be run on the CPU or GPU, with no additional licenses required.

**Participated benchmarks** `ACASXu`, `cifar10_resnet`, `cifar2020`, `eran`, `marabou-cifar10`, `mnistfc`, `nn4sys`, `oval21`, `verivital`.

## 3.6  NNV

**Team** Neelanjana Pal (Vanderbilt University), Taylor T Johnson (Vanderbilt University)

**Description** The Neural Network Verification Tool (NNV) [43–46, 52] is written primarily with Matlab and implements reachability-analysis methods for neural network verification with a particular focus on applications of closed-loop neural network control systems in autonomous cyber-physical systems. NNV uses geometric representations such as star sets that allows for a layer-by-layer computation of the exact reachable set for feed-forward deep neural networks. In the event that a particular safety property is violated, NNV can construct and visualize the complete set of counterexample inputs for a neural network.

**Link** https://github.com/verivital/nnv.

**Commit** `3ca2629aaceb9080e4d08a0f9c6b51854f9c7b7f` (for reproducing competition results; otherwise please use the master version).

**Hardware and licences** GPU and CPU with 64-bit floating point. A license for Matlab will be required.

**Participated benchmarks** `ACASXu`, `cifar2020`, `eran`, `mnistfc`, `oval21`, `verivital`.

## 3.7  OVAL

**Team** Alessandro De Palma (University of Oxford), Florian Jaeckle (University of Oxford), M. Pawan Kumar (University of Oxford)

**Description** The OVAL verification tool is an optimization-based complete verifier that relies on a specialized Branch and Bound (BaB) framework for neural network verification. In this context, a BaB method is composed of three main components (see [9] for an overview): a *branching strategy* to divide the verification property into easier subproblems, a *bounding algorithm* to compute over-approximation bounds for each subproblem, and a *falsification algorithm* to look for counter-examples to the property.

For networks with medium to large input dimensionality, OVAL relies on the efficient FSB [13] branching strategy, which combines a dual-based scoring of ReLU neurons [9] with inexpensive strong branching approximations to select an activation to split upon. For networks of small input dimensionality, OVAL can revert to input splitting as in [10]. Bounds on the

subproblems are obtained by adaptively choosing [12] between bounding algorithms of varying degrees of tightness: the competition entry relies on Beta-CROWN [49], which effectively solves the convex hull of element-wise activations, and Active Set [11], which operates on the tighter relaxation from [2] to tackle harder properties. In addition, the framework supports a variety of bounding algorithms [8, 12, 15, 54]. The search for counter-examples is performed using the MI-FGSM [14] adversarial attack, which we adapted to perform general property falsification.

The implementation of the OVAL framework, written in PyTorch [31], exploits GPU acceleration and is massively parallel over both the BaB subproblems and the relative intermediate computations [8]. It currently supports fully connected and convolutional networks, with ReLU, maxpool and average pooling layers.

**Link** https://github.com/oval-group/oval-bab.

**Commit** `014b6ee5071508430c8e515bbae725306db68fe1` in order to reproduce competition results. We otherwise suggest to employ the master version.

**Hardware and licences** GPU and CPU with 32-bit or 64-bit floating point. No license is required.

**Participated benchmarks** ACASXu, `cifar2020`, `eran`, `marabou-cifar10`, `mnistfc`, `nn4sys`, `oval21`, `verivital`.

## 3.8 RPM

**Team** Joe Vincent (Stanford), Mac Schwager (Stanford)

**Description** The Reachable Polyhedral Marching (RPM) tool [47] is a method for computing exact forward and backward reachable sets of feedforward neural networks with ReLU activation. Verification problems are posed as backward reachability problems. A unique feature of the RPM tool is its incremental computation of reachable sets. For verification this means that unsafe inputs may be found before computing the complete reachable set, leading to early termination. RPM does not currently have a parallel implementation, although the algorithm is amenable to parallelization.

**Link** https://github.com/StanfordMSL/Neural-Network-Reach/tree/vnn_comp_2021

**Commit** `861ce6e380e3cc2d439a7bca87b59817e4624af6` for reproducing competition results. For other purposes the most recent commit is suggested.

**Hardware and licences** CPU, no license is required.

**Participated benchmarks** ACASXu

## 3.9 ComposableNeuralVerification (NV.jl)

**Team** Tianhao Wei (Carnegie Mellon), Chen Tan (Northeastern), Changliu Liu (Carnegie Mellon)

**Description** This tool is adapted from the original NeuralVerification.jl [25] developed at the Stanford Intelligent System Lab. This Julia toolbox implemented a wide variety of verification algorithms that use reachability, optimization, and search, which are summarized in [26]. We added support for onnx format networks and vnnlib format specifications,

**Link** https://github.com/intelligent-control-lab/NeuralVerification.jl

**Hardware and licences** CPU, no licence

**Commit** `4e612602ba4b34b42416742d85476d9b0dcdcb51` (for reproducing competition results; otherwise please use the master branch)

**Participated benchmarks** nn4sys and AcasXu

## 3.10 Venus

**Team** Panagiotis Kouvaros (Imperial College London), Alessio Lomuscio (Imperial College London)

**Description** Venus is a complete verification tool for Relu-based feed-forward neural networks. Venus implements a MILP-based verification method whereby it leverages dependency relations between the ReLU nodes to reduce the search space that needs to be considered during branch-and-bound. The dependency relations are exploited via callback cuts [6] and via a branching method that divides the verification problem into a set of sub-problems whose MILP formulations require fewer integrality constraints [23]. To derive tight MILP encodings, Venus additionally implements a symbolic interval propagation method for computing the pre-activation bounds of the ReLU nodes; the method optimises the linear relaxation of each of the ReLU nodes towards minimising the over-approximation error in subsequent layers.

**Link** https://github.com/vas-group-imperial/venus2

**Commit** For the reproduction of the VNN-COMP2021 results please use the repository https://github.com/pkouvaros/venus2_vnncomp21 (57e9608041d230b5d78c4f2afb890b81035436a1).

**Hardware and licenses** CPU, GUROBI License.

**Participated benchmarks** ACASXU, mnistfc, nn4sys.

## 3.11 Debona

**Team** Christopher Brix (RWTH Aachen University), Thomas Noll (RWTH Aachen University)

**Description** Debona is a fork of VeriNet [17]. However, the abstract domain used by VeriNet defines symbolic linear upper and lower bounds that are parallel to each other, i.e., offset only by some scalar value. On the contrary, Debona utilizes independent upper and lower bounds. This allows for a tighter relaxation especially for ReLU operations, where a lower bound of zero may be better than bounds that are negative in large regions of the input space. This idea has been described in [7] but was independently previously published in [38].

**Link** https://github.com/ChristopherBrix/Debona

**Commit** f000f3d483b2cc592233d0ba2a1a0327210562c8

**Hardware and licences** CPU, Gurobi licence

**Participated benchmarks** AcasXu, eran, mnistfc and nn4sys

## 3.12 nnenum

**Team** Stanley Bak (Stony Brook Univeristy)

**Description** The nnenum tool uses multiple levels of abstraction to achieve high-performance verification of ReLU networks without sacrificing completeness [3]. Analysis combines three types of zonotopes with star set (triangle) overapproximations [45], and uses efficient parallelized ReLU case splitting [5]. The ImageStar method [43] allows sets to be quickly propagated through all layers supported by the ONNX runtime, such as convolutional layers with arbitrary parameters. The tool is written in Python 3, uses GLPK for LP solving. New this year we added support for `vnnlib` files, and optimized some of the LP timeout parameters for acasxu [4].

**Link** https://github.com/stanleybak/nnenum

**Commit** c93a39cb568f58a26015bd151acafab34d2d4929

| Benchmark | Application | Network Types | Largest NN |
|---|---|---|---|
| Acasxu | Control | Feedforward + ReLU Only | 54.6k |
| Cifar10_resnet | Image Classification | ResNet | 487k |
| Cifar2020 (unscored) | Image Classification | Conv + ReLU | 9.41M |
| Eran | Image Classification | Feedforward + non-ReLU | 1.68M |
| Marabou-cifar10 | Image Classification | Conv + ReLU | 1.29M |
| Mnistfc | Image Classification | Feedforward + ReLU Only | 2.03M |
| nn4sys | Database Indexing | Feedforward + ReLU Only | 336.5M* |
| Oval21 | Image Classification | Conv + ReLU | 840k |
| Verivital | Image Classification | Conv + maxpool / avgpool | 46.3k |

*After zipping, the network is of 1.79M.

Table 2: Overview of all benchmarks.

**Hardware and licences** CPU, No licenses required

**Participated benchmarks** AcasXu, cifar2020, mnistfc, oval

# 4 Benchmarks

## 4.1 ACAS Xu

**Networks** The ACASXu benchmarks consists of ten properties defined over 45 neural networks used to issue turn advisories to aircraft to avoid collisions. The neural networks have 300 neurons arranged in 6 layers, with ReLU activation functions. There are five inputs corresponding to the aircraft states, and five network outputs, where the minimum output is used as the turn advisory the system ultimately produces.

**Specifications** We use the original 10 properties [21], where properties 1-4 are checked on all 45 networks as was done in later work by the original authors [22]. Properties 5-10 are checked on a single network. The total number of benchmarks is therefore 186. The original verification times ranged from seconds to days—including some benchmark instances that did not finish. This year we used a timeout of around two minutes (116 seconds) for each property, in order to fit within a total maximum runtime of six hours.

## 4.2 Cifar10_resnet

**Proposed by** the $\alpha,\beta$-CROWN team.

**Motivations** Currently, many tools are hard-coded to handle feedforward networks only. To make neural network verification more useful in practical scenarios, we advocate that tools should handle more general architecture. Residual networks [16] (ResNet) is the first step towards this goal due to its relatively simple structure and practical significance. The propose of this benchmark is to provide some incentives for the community to develop more generic tools.

**Networks** We provided two ResNet models on CIFAR-10 image classification task with the following structures:

- `ResNet-2B` with 2 residual blocks: 5 convolutional layers + 2 linear layers

- `ResNet-4B` with 4 residual blocks: 9 convolutional layers + 2 linear layers

The networks are trained via adversarial training with an $\ell_\infty$ perturbation norm of $\epsilon = \frac{2}{255}$. For simplicity, these networks do not contain batch normalization or pooling layers, and use ReLU activation functions. The ResNet-4B model is relatively large with over 10K neurons.

We evaluated both networks using a 100-step projected gradient descent (PGD) attack with 5 random restarts, and a simple bound propagation based verification algorithm CROWN [55] (mathematically equivalent to the abstract interpretation used in DeepPoly [38]) under $\ell_\infty$ norm perturbations. The results are listed in Table 3.

| Model | # ReLUs | Clean acc. | $\epsilon = 2/255$ | | $\epsilon = 1/255$ | |
| --- | --- | --- | --- | --- | --- | --- |
| | | | PGD acc. | Verified acc. | PGD acc. | Verified acc. |
| `ResNet-2B` | 6244 | 69.25% | 54.82% | 26.88% | 62.24% | 57.16% |
| `ResNet-4B` | 14436 | 77.20% | 61.41% | 0.24% | 69.75% | 23.28% |

Table 3: Clean accuracy, PGD accuracy and CROWN verified accuracy for ResNet models. Note that the verified accuracy is obtained via the vanilla version of CROWN/DeepPoly which has been widely used as a simple baseline, not the $\alpha,\beta$-CROWN tool used in the competition.

To ensure the appropriate level of difficulty, we use $\epsilon = \frac{2}{255}$ for the ResNet-2B model and $\epsilon = \frac{1}{255}$ for the ResNet-4B model.

**Specifications** We randomly select 48 images from the CIFAR-10 test set for the ResNet-2B model and 24 images for the ResNet-4B model. The images are classified correctly and cannot be attacked by a 100-step PGD attack with 5 random restarts. For each image, we specify the property that the logit of the ground-truth label is always greater than the logits of all other 9 labels within $\ell_\infty$ norm input perturbation of $\epsilon = \frac{2}{255}$ for ResNet-2B and $\epsilon = \frac{1}{255}$ for ResNet-4B. The per-example timeout is set to 5 minutes and the overall runtime is guaranteed to be less than 6 hours.

## 4.3 Cifar2020 (unscored)

**Motivation** This benchmark combines two convolutional CIFAR10 networks from last year's VNN-COMP 2020 with a new, larger network with the goal to evaluate the progress made by the whole field of Neural Network verification.

**Networks** The two ReLU networks `cifar_10_2_255` and `cifar_10_8_255` with two convolutional and two fully-connected layers were trained for $\ell_\infty$ perturbations of $\epsilon = \frac{2}{255}$ and $\frac{8}{255}$, respectively, using COLT [29] and the larger `ConvBig` with four convolutional and three fully-connected networks, was trained using adversarial training [28] and $\epsilon = \frac{2}{255}$.

**Specifications** We draw the first 100 images from the CIFAR10 test set and for every network reject incorrectly classified ones. For the remaining images, the specifications describe a correct classification under an $\ell_\infty$-norm perturbation of at most $\frac{2}{255}$ and $\frac{8}{255}$ for `cifar_10_2_255` and `ConvBig` and `cifar_10_8_255`, respectively and allow a per sample timeout of 5 minutes.

## 4.4 eran

**Proposed by:** The ERAN team

**Motivation** While most Neural Network Verificaiton methods focus their analysis on ReLU based networks, many modern network architectures, e.g., EfficientNet [41], are based on non-piecwise-linear activation functions. To begin to understand how the choice of activation function affects certifiability, the eran benchmark aims at comparing the certifiabilty of networks based on piecewise-linear and non-piecewise-linear activation functions under an $\ell_\infty$-norm based adversary.

**Networks** We consider a ReLU network with 8 hidden layers of width 200 and a Sigmoid network with 6 hidden layers of width 200. Both networks were trained using standard training.

**Specifications** We sample random images from the MNIST test set until we obtain 36 correctly classified images per network. For these images, the specifications describe a correct classification under an $\ell_\infty$-norm perturbation of at most 0.015 and 0.012 for the ReLU and Sigmoid network, respectively, and allow a per sample timeout of 5 minutes.

## 4.5   Marabou-cifar10

**Proposed by** The Marabou team.

**Networks** This benchmark contains three convolutional networks, `cifar10_small.onnx`, `cifar10_medium.onnx`, and `cifar10_large.onnx`, trained on the CIFAR10 dataset. Each network has 2 convolutional layers followed by two fully connected feed-forward layers. Each layer uses the ReLU activation functions. The networks are all trained with Adam optimizer for 120 epochs with learning rate 0.0002. The three networks contain 2568, 4944, and 10528 ReLUs, respectively. The test accuracy are 63.14%, 70.21%, and 74.16%, respectively. The networks expect the input image to be normalized between 0 and 1.

**Specifications** We randomly sample correctly classified images from the CIFAR10 test set. The specifications are targeted adversarial robustness, which states that the network does not mis-classify an image as a given adversarial label under $l_\infty$-norm perturbations. The target label is chosen as $(correctLabel + 1) \mod 10$. We propose two perturbation bounds: 0.012 and 0.024, and allow a per-query timeout of 5 minutes.

## 4.6   Mnistfc

**Proposed by** The VeriNet team.

**Motivation** This benchmark contains fully connected networks with ReLU activation functions and varying depths.

**Networks** The benchmark set consists of three fully-connected classification networks with 2, 4 and 6 layers and 256 ReLU nodes in each layer trained on the MNIST dataset. The networks were first presented in a benchmark in VNN-COMP 2020 [48].

**Specifications** We randomly sampled 15 correctly classified images from the MNIST test set. For each network and image, the specification was a correct classification under $l_\infty$ perturbations of at most $\epsilon = 0.03$ and $\epsilon = 0.05$. The timeouts were 2 minutes per instance for the 2-layer network and 5 minutes for the remaining two networks.

## 4.7   NN4Sys

**Proposed by** The ComposableNeuralverification team

**Application** The benchmark contains networks for database indexing, which is a 1D to 1D mapping.

- *Background*: learned index is a neural network (NN) based database index proposed by Kraska et al. [24], 2018. It shows great potential but has one drawback—for non-existing keys (i.e., the keys that do not exist in the database), the outputs of a learned index can be arbitrary.

- *What we do*: to provide safety guarantees for *all* keys, we design a specification to dictate how "far" one predicted position can be, compared to its actual position (or the positions that non-existing keys should be).

- *What to verify*: our benchmark provides multiple pairs of (1) learned indexes (trained NNs) and (2) corresponding specifications. We design these pairs with different parameters such that they cover a variety of user needs and have varied difficulties for verifiers.

- *Translating learned indexes to a VNN benchmark*: the original learned index [24] contains two stages (of NNs) for high precision. However, this cascading structure is inconvenient/unsupported to verify because there is a "switch" operation—choosing one NN in the second stage based on the prediction of the first stage's NN. To convert learned indexes to a standard form, we merge the NNs in both stages into an integrated network by adding some hand crafted layers.

- *A note on broader impact*: using NNs for systems is a broad topic, but many existing works lack strict safety guarantees. We believe that NN Verification can help system developers gain confidence to apply NNs to critical systems. We hope our benchmark can be an early step towards this vision.

**Networks** The networks are feedforward with ReLU activations, and they are sparse networks. There are six networks in this benchmark. Three of them have original size 194.2M and zipped size 1.79M; the other three have original size 336.5M and zipped size 790k. This is because onnx does not support directly encoding of sparse matrices, hence the networks are stored as fully connected networks.

**Specifications** The specification aims to check if the prediction error is bounded. The specification is a collection of pairs of input and output intervals such that any input in the input interval should be mapped to the corresponding output interval.

## 4.8 Oval21

**Proposed by** The OVAL team.

**Motivations** The majority of adversarial robustness benchmarks consider image-independent perturbation radii, possibly resulting in some properties that are either easily verified by all verification methods, or too hard to be verified (for commonly employed timeouts) by any of them. In line with the OVAL verification dataset from VNN-COMP 2020 [48], whose versions have already been used in various recent works [8, 11–13, 19, 20, 27, 49, 54], the OVAL21 benchmark associates to each image-network pair a perturbation radius found via binary search to ensure that all properties are challenging to solve.

**Networks** The benchmark includes 3 ReLU-based convolutional networks which were robustly trained [50] against $\ell_\infty$ perturbations of radius $\epsilon = 2/255$ on CIFAR10. Two of the networks, named `base` and `wide`, are composed of 2 convolutional layers followed by 2 fully connected layers and have respectively 3172 and 6244 activations. The third model, named `deep`, has 2 additional convolutional layers and a total of 6756 activations.

**Specifications** The verification properties represent untargeted adversarial robustness (with respect to all possible misclassifications) to $\ell_\infty$ perturbations of varying $\epsilon$, with a per-instance timeout of 720 seconds. The property generation procedure relies on commonly employed lower and upper bounds to the adversarial loss to exclude perturbation radii that yield trivial properties. 10 correctly classified images per network are randomly sampled from the entire CIFAR10 test set, and a distinct $\epsilon \in [0, 16/255]$ is associated to each. First, a binary search is run to find the largest $\epsilon$ value for which a popular iterative adversarial attack [14] fails to find an adversarial example. Then, a second binary search is run to find the smallest $\epsilon$ value for which bounds [54] from the element-wise convex hull of the activations (with fixed intermediate bounds from [50, 55]) fail to prove robustness. Both binary search procedures are run with a tolerance of $\epsilon_{\text{tol}} = 0.1$. Denoting $\epsilon_{lb}$ as the smallest output from the two routines, and $\epsilon_{ub}$ as the largest, the following perturbation radius is chosen: $\epsilon = \frac{1}{3}\epsilon_{lb} + \frac{2}{3}\epsilon_{ub}$.

**Link** https://github.com/stanleybak/vnncomp2021/tree/main/benchmarks/oval21

## 4.9 Verivital

**Proposed by** The VeriVITAL team.

**Motivation** Neural networks with pooling layers are vastly used in several applications. The main motivation for proposing this benchmark was to include the pooling layers as part of this year's VNN-Comp.

**Networks** This benchmark contains two MNIST classifiers with pooling layers, one with averagepooling layers and the other with maxpooling.

**Specifications** We randomly sampled 20 correctly classified images from the MNIST test set. For the network with averagepooling layers, the specification was to correctly classify those randomly chosen images with an $l_\infty$ perturbation radii($\epsilon$) of 0.02 and 0.04 and a timeout of 5 minutes. For the network with maxpooling layers the corresponding radius was 0.004 with a timeout of 7 minutes.

**Link** https://github.com/stanleybak/vnncomp2021/tree/main/benchmarks/verivital

# 5   Results

Each tool was run on all benchmarks which produced a `csv` file of results. This was sent to the authors for review, which sometimes required rerunning certain benchmarks to make sure they match the expected results. Python code was then written to process the results and compute the scores. The final `csv` files for each tool as well as scoring scripts are available online: https://github.com/stanleybak/vnncomp2021_results.

This also includes detailed log files for each benchmark showing the specific runtime for each tool and the score awarded. This can be used to find challenging instances to help with tool development. For example, in the output files in the repo you may see things like:

```
Row: ['ACASXU_run2a_4_2_batch_2000-prop_2', '-', '6.4 (h)', '10.5 (h)',
      'timeout', '41.1 (h)', 'timeout', 'timeout', '64.8 (h)', '62.5 (h)',
      'timeout', 'timeout', 'timeout', '-']
73: nnv score: 0
73: nnenum score: 12
73: venus2 score: 11
73: NN-R score: 0
73: VeriNet score: 10
73: DNNF score: 0
73: Debona score: 0
73: a-b-CROWN score: 10
73: oval score: 10
73: Marabou score: 0
73: ERAN score: 0
73: NV.jl score: 0
73: randgen score: 0
```

The tools are listed in order, and row is the times and result for each tool. So `nnenum` should be holds with a time of 6.4 (after subtracting overhead). The scores are also listed for each tool. Since `nnenum` was the fastest on this instance, it got 12 points, 10 for correct plus 2 for time bonus as fastest. The second fastest was `venus2` at 10.5 seconds, so they get 11 points. None of the remaining tools were within 0.2 seconds, so they all received 10 points if they completed analysis successfully.

## 5.1   Overall Score

The overall score for VNN-COMP 2021 is shown in Table 4. Two tables are included, based on the two ways to measure overhead discussed in Section 2. Overall, $\alpha,\beta$-CROWN performed best, followed by VeriNet. We awarded two third place results, one to oval and one to ERAN, as their ranking depended upon factors like overhead as well as how incorrect results were judged, discussed next. For all the remaining tables in this section, the numbers reported correspond to the multi-overhead measurements.

One unexpected aspect was how to judge incorrect results, since tools currently are not required to produce a counter-example when an instance is falsified. We considered two reasonable options, which was voting (majority is assumed to be correct), and odd-one-out. In odd-one-out, only if a single tool's output differs from all the others is the result is considered incorrect. If multiple tools produce the same result or if only two tools completed the instance and their results differ, then the instance is ignored for scoring purposes. This generally had a

slight effect on the score, which was significant enough to affect the order in the total ranking. Specifically, the positions of ERAN and oval for third place could be affected by the scoring parameters, when using the "Single Overhead" overhead correction. Notice that, however, both voting and odd-one-out are imperfect ways to judge incorrect results, and it may be the case that the mismatching tool was in fact correct. In future editions of VNN-COMP we may standardize counter-example outputs and require they are produced when instances can be falsified, to remedy this shortcoming. For the reported category scores, we display odd-one-out scores and results.

Other statistics and the individual benchmark scores are also included below. In general, the GPU tools did better, as well as tools that could support a large number of benchmarks and verified a large number of benchmark instances. Several tools produced mismatching results, and an interesting followup study could identify the underlying reasons for this unsoundness.

Table 4: VNN-COMP 2021 Overall Score

(a) Single Overhead

| # | Tool | Score |
|---|---|---|
| 1 | $\alpha,\beta$-CROWN | 779.2 |
| 2 | VeriNet | 701.2 |
| 3 | oval | 582.0 |
| 4 | ERAN | 581.1 |
| 5 | Marabou | 335.3 |
| 6 | Debona | 201.9 |
| 7 | venus2 | 189.2 |
| 8 | nnenum | 184.5 |
| 9 | nnv | 57.2 |
| 10 | NV.jl | 48.1 |
| 11 | RPM | 25.4 |
| 12 | DNNF | 24.3 |
| 13 | randgen | 1.9 |

(b) Multi-Overhead

| # | Tool | Score |
|---|---|---|
| 1 | $\alpha,\beta$-CROWN | 779.7 |
| 2 | VeriNet | 705.0 |
| 3 | ERAN | 643.4 |
| 4 | oval | 581.8 |
| 5 | Marabou | 339.0 |
| 6 | Debona | 201.9 |
| 7 | venus2 | 189.2 |
| 8 | nnenum | 184.6 |
| 9 | nnv | 57.2 |
| 10 | NV.jl | 48.1 |
| 11 | RPM | 25.4 |
| 12 | DNNF | 24.3 |
| 13 | randgen | 1.9 |

## 5.2   Other Stats

This section presents other statistics related to the measurements that are interesting, but did not play a direct role in scoring this year. With ERAN, which had different overheads, the 'CPU Mode' column in Table 5 corresponds to the overhead used for the ACASXu and ERAN benchmarks, whereas the 'Seconds' column corresponds to when the GPU was used (all others).

Table 5: Overhead

| # | Tool | Seconds | CPU Mode |
|---|---|---|---|
| 1 | Marabou | 0.2 | - |
| 2 | randgen | 0.3 | - |
| 3 | nnenum | 1.0 | - |
| 4 | venus2 | 1.7 | - |
| 5 | DNNF | 2.0 | - |
| 6 | VeriNet | 2.2 | - |
| 7 | Debona | 2.5 | - |
| 8 | oval | 5.1 | - |
| 9 | $\alpha,\beta$-CROWN | 6.1 | - |
| 10 | ERAN | 7.1 | 3.7 |
| 11 | nnv | 8.4 | - |
| 12 | NV.jl | 20.9 | - |
| 13 | RPM | 52.2 | - |

Table 6: Num Benchmarks Participated

| # | Tool | Count |
|---|---|---|
| 1 | VeriNet | 9 |
| 2 | ERAN | 9 |
| 3 | $\alpha,\beta$-CROWN | 9 |
| 4 | oval | 8 |
| 5 | Marabou | 7 |
| 6 | nnv | 6 |
| 7 | DNNF | 5 |
| 8 | randgen | 4 |
| 9 | nnenum | 4 |
| 10 | Debona | 4 |
| 11 | venus2 | 3 |
| 12 | NV.jl | 2 |
| 13 | RPM | 1 |

Table 7: Num Instances Verified

| # | Tool | Count |
|---|------|-------|
| 1 | $\alpha,\beta$-CROWN | 766 |
| 2 | VeriNet | 717 |
| 3 | ERAN | 656 |
| 4 | oval | 636 |
| 5 | Marabou | 364 |
| 6 | nnenum | 310 |
| 7 | Debona | 280 |
| 8 | venus2 | 266 |
| 9 | nnv | 141 |
| 10 | NV.jl | 97 |
| 11 | RPM | 72 |
| 12 | DNNF | 55 |
| 13 | randgen | 33 |

Table 8: Num Violated

| # | Tool | Count |
|---|------|-------|
| 1 | ERAN | 177 |
| 2 | $\alpha,\beta$-CROWN | 177 |
| 3 | oval | 175 |
| 4 | VeriNet | 175 |
| 5 | Marabou | 103 |
| 6 | nnenum | 73 |
| 7 | Debona | 70 |
| 8 | venus2 | 63 |
| 9 | DNNF | 55 |
| 10 | RPM | 44 |
| 11 | randgen | 33 |
| 12 | NV.jl | 11 |

Table 9: Num Holds

| # | Tool | Count |
|---|------|-------|
| 1 | $\alpha,\beta$-CROWN | 589 |
| 2 | VeriNet | 542 |
| 3 | ERAN | 479 |
| 4 | oval | 461 |
| 5 | Marabou | 261 |
| 6 | nnenum | 237 |
| 7 | Debona | 210 |
| 8 | venus2 | 203 |
| 9 | nnv | 141 |
| 10 | NV.jl | 86 |
| 11 | RPM | 28 |

Table 10: Mismatched (Incorrect) Results

| # | Tool | Count |
|---|------|-------|
| 1 | Marabou | 26 |
| 2 | Debona | 14 |
| 3 | nnv | 12 |
| 4 | NV.jl | 5 |
| 5 | venus2 | 1 |

## 5.3 Benchmark Scores

The results for the individual categories are shown below. For overall score, the tools which participated in all or almost all of the benchmarks did best. Within individual benchmarks, some tools performed well despite not ranking high in the overall score. The results presented here are for multi-overhead setup, with incorrect results scored using the odd-one-out strategy. Adjusting these parameters produced minor changes in the overall score and rankings, but was omitted for clarity. If these alternate scores are of interest, the discussion at the beginning of Section 5 outlines where to access the scripts used to compute scores.

Table 11: Benchmark `acasxu`

| # | Tool | Verified | Falsified | Fastest | Score | Percent |
|---|------|----------|-----------|---------|-------|---------|
| 1 | nnenum | 138 | 47 | 155 | 1910 | 100.0% |
| 2 | VeriNet | 138 | 47 | 117 | 1852 | 97.0% |
| 3 | Marabou | 137 | 46 | 115 | 1809 | 94.7% |
| 4 | oval | 138 | 47 | 98 | 1794 | 93.9% |
| 5 | venus2 | 138 | 46 | 94 | 1778 | 93.1% |
| 6 | $\alpha,\beta$-CROWN | 138 | 47 | 67 | 1732 | 90.7% |
| 7 | ERAN | 125 | 46 | 24 | 1506 | 78.8% |
| 8 | Debona | 84 | 42 | 39 | 1086 | 56.9% |
| 9 | RPM | 28 | 44 | 9 | 486 | 25.4% |
| 10 | nnv | 29 | 0 | 29 | 348 | 18.2% |
| 11 | DNNF | 0 | 41 | 12 | 182 | 9.5% |
| 12 | randgen | 0 | 28 | 0 | 28 | 1.5% |
| 13 | NV.jl | 45 | 9 | 0 | -23 | 0% |

Table 12: Benchmark `cifar10-resnet`

| # | Tool | Verified | Falsified | Fastest | Score | Percent |
|---|------|----------|-----------|---------|-------|---------|
| 1 | $\alpha,\beta$-CROWN | 58 | 0 | 12 | 623 | 100.0% |
| 2 | VeriNet | 48 | 0 | 29 | 548 | 88.0% |
| 3 | ERAN | 43 | 0 | 36 | 502 | 80.6% |
| 4 | Marabou | 39 | 0 | 0 | 390 | 62.6% |

Table 13: Benchmark `cifar2020`

| # | Tool | Verified | Falsified | Fastest | Score | Percent |
|---|------|----------|-----------|---------|-------|---------|
| 1 | oval | 146 | 41 | 174 | 2209 | 100.0% |
| 2 | $\alpha,\beta$-CROWN | 148 | 42 | 43 | 1996 | 90.4% |
| 3 | VeriNet | 139 | 42 | 5 | 1822 | 82.5% |
| 4 | ERAN | 107 | 43 | 132 | 1749 | 79.2% |
| 5 | nnenum | 62 | 13 | 0 | 741 | 33.5% |
| 6 | randgen | 0 | 2 | 0 | 2 | 0.1% |
| 7 | nnv | 6 | 0 | 0 | -140 | 0% |

Table 14: Benchmark `eran`

| # | Tool | Verified | Falsified | Fastest | Score | Percent |
|---|------|----------|-----------|---------|-------|---------|
| 1 | $\alpha,\beta$-CROWN | 60 | 1 | 26 | 670 | 100.0% |
| 2 | VeriNet | 48 | 1 | 49 | 588 | 87.8% |
| 3 | ERAN | 46 | 1 | 0 | 470 | 70.1% |
| 4 | Debona | 47 | 2 | 39 | 375 | 56.0% |
| 5 | oval | 21 | 0 | 18 | 247 | 36.9% |
| 6 | Marabou | 19 | 0 | 0 | 190 | 28.4% |
| 7 | nnv | 10 | 0 | 0 | 100 | 14.9% |

Table 15: Benchmark `marabou-cifar10`

| # | Tool | Verified | Falsified | Fastest | Score | Percent |
|---|------|----------|-----------|---------|-------|---------|
| 1 | $\alpha,\beta$-CROWN | 1 | 52 | 52 | 625 | 100.0% |
| 2 | ERAN | 0 | 52 | 51 | 613 | 98.1% |
| 3 | oval | 0 | 53 | 44 | 611 | 97.8% |
| 4 | VeriNet | 0 | 52 | 16 | 543 | 86.9% |
| 5 | Marabou | 0 | 29 | 0 | 281 | 45.0% |
| 6 | DNNF | 0 | 2 | 0 | 11 | 1.8% |
| 7 | randgen | 0 | 1 | 0 | 1 | 0.2% |

Table 16: Benchmark `mnistfc`

| # | Tool | Verified | Falsified | Fastest | Score | Percent |
|---|------|----------|-----------|---------|-------|---------|
| 1 | $\alpha,\beta$-CROWN | 49 | 21 | 35 | 772 | 100.0% |
| 2 | VeriNet | 39 | 21 | 57 | 716 | 92.7% |
| 3 | Debona | 37 | 22 | 47 | 688 | 89.1% |
| 4 | oval | 37 | 21 | 46 | 676 | 87.6% |
| 5 | ERAN | 34 | 22 | 47 | 654 | 84.7% |
| 6 | Marabou | 35 | 19 | 0 | 540 | 69.9% |
| 7 | venus2 | 31 | 16 | 26 | 522 | 67.6% |
| 8 | nnenum | 35 | 12 | 21 | 512 | 66.3% |
| 9 | nnv | 27 | 0 | 8 | 186 | 24.1% |
| 10 | DNNF | 0 | 7 | 0 | 70 | 9.1% |

Table 17: Benchmark `nn4sys`

| # | Tool | Verified | Falsified | Fastest | Score | Percent |
|---|------|----------|-----------|---------|-------|---------|
| 1 | $\alpha,\beta$-CROWN | 70 | 5 | 73 | 878 | 100.0% |
| 2 | VeriNet | 68 | 3 | 0 | 719 | 81.9% |
| 3 | ERAN | 67 | 4 | 0 | 704 | 80.2% |
| 4 | oval | 56 | 4 | 0 | 589 | 67.1% |
| 5 | NV.jl | 41 | 2 | 0 | 422 | 48.1% |
| 6 | venus2 | 34 | 1 | 0 | 250 | 28.5% |
| 7 | DNNF | 0 | 4 | 0 | 22 | 2.5% |
| 8 | randgen | 0 | 2 | 0 | 2 | 0.2% |
| 9 | Debona | 42 | 4 | 0 | -722 | 0% |

Table 18: Benchmark `oval21`

| # | Tool | Verified | Falsified | Fastest | Score | Percent |
|---|------|----------|-----------|---------|-------|---------|
| 1 | oval | 12 | 2 | 11 | 164 | 100.0% |
| 2 | $\alpha,\beta$-CROWN | 12 | 2 | 2 | 146 | 89.0% |
| 3 | VeriNet | 11 | 2 | 6 | 145 | 88.4% |
| 4 | ERAN | 6 | 2 | 3 | 86 | 52.4% |
| 5 | Marabou | 4 | 2 | 1 | 63 | 38.4% |
| 6 | nnenum | 2 | 1 | 0 | 30 | 18.3% |
| 7 | nnv | 16 | 0 | 4 | -31 | 0% |

Table 19: Benchmark `verivital`

| # | Tool | Verified | Falsified | Fastest | Score | Percent |
|---|------|----------|-----------|---------|-------|---------|
| 1 | $\alpha,\beta$-CROWN | 53 | 7 | 52 | 704 | 100.0% |
| 2 | oval | 51 | 7 | 57 | 694 | 98.6% |
| 3 | ERAN | 51 | 7 | 56 | 693 | 98.4% |
| 4 | VeriNet | 51 | 7 | 0 | 580 | 82.4% |
| 5 | DNNF | 0 | 1 | 0 | 10 | 1.4% |
| 6 | nnv | 53 | 0 | 0 | -168 | 0% |
| 7 | Marabou | 27 | 7 | 0 | -2260 | 0% |

# 6 Conclusion and Ideas for Future Competitions

This report summarizes the 2ⁿᵈ Verification of Neural Networks Competition (VNN-COMP) held in 2021. Improvements to the competition structure have been made, including standardization of common input formats (`onnx` and `vnnlib`), and common measurement hardware. Based on the common benchmarks (CIFAR2020 and ACASXU), tools exhibited significant progress in terms of scalability and speed compared with previous years. The comparison is imperfect, as last year we did not have standardized hardware, so it is unclear how much of the speed improvement is due to algorithmic improvements and how much is due to better hardware. Future editions of the competition may better judge the year-to-year improvements in neural network verification methods.

The benchmarks and tool execution scripts are openly available for others to replicate: https://github.com/stanleybak/vnncomp2021. We hope this serves as a fair comparison for evaluating future improvements to verification methods in upcoming publications. From an applicability perspective, having a common input format hopefully reduces the barriers of industry participants to use the developed tools.

In future editions of VNN-COMP, some improvements we can make would be to allow more types of hardware. For example, any AWS EC2 machine could be chosen by the tool authors with roughly the same cost ($3 an hour this year). Multiple tool authors expressed that their tool could work faster if given both a strong CPU and GPU together. Alternatively, we could allow custom hardware per benchmark. This would likely require additional automation with the competition measurements, especially dealing with installing license files. Some ideas for automating this could be to have authors provide the Gurobi licences to use, rather than the competition organizers. Another improvement would be to improve overhead measurement, as detailed in Section 2. Finally, we should standardize the counter-example format, so that the ground truth for mismatched verification results can be provided.

# Acknowledgements

Tool authors listed in section 3 participated in the preparation and review of this report.

# References

[1] Guy Amir, Haoze Wu, Clark Barrett, and Guy Katz. An smt-based approach for verifying binarized neural networks. *arXiv preprint arXiv:2011.02948*, 2020.

[2] Ross Anderson, Joey Huchette, Will Ma, Christian Tjandraatmadja, and Juan Pablo Vielma. Strong mixed-integer programming formulations for trained neural networks. *Mathematical Programming*, 2020.

[3] Stanley Bak. Execution-guided overapproximation (ego) for improving scalability of neural network verification, 2020.

[4] Stanley Bak. nnenum: Verification of relu neural networks with optimized abstraction refinement. In *NASA Formal Methods Symposium*, pages 19–36. Springer, 2021.

[5] Stanley Bak, Hoang-Dung Tran, Kerianne Hobbs, and Taylor T. Johnson. Improved geometric path enumeration for verifying ReLU neural networks. In *32nd International Conference on Computer-Aided Verification (CAV)*, July 2020.

[6] E. Botoeva, P. Kouvaros, J. Kronqvist, A. Lomuscio, and R. Misener. Efficient verification of neural networks via dependency analysis. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence (AAAI20)*. AAAI Press, 2020.

[7] Christopher Brix and Thomas Noll. Debona: Decoupled boundary network analysis for tighter bounds and faster adversarial robustness proofs. *CoRR*, abs/2006.09040, 2020.

[8] Rudy Bunel, Alessandro De Palma, Alban Desmaison, Krishnamurthy Dvijotham, Pushmeet Kohli, Philip HS Torr, and M Pawan Kumar. Lagrangian decomposition for neural network verification. *Conference on Uncertainty in Artificial Intelligence*, 2020.

[9] Rudy Bunel, Jingyue Lu, Ilker Turkaslan, P Kohli, P Torr, and M Pawan Kumar. Branch and bound for piecewise linear neural network verification. *Journal of Machine Learning Research*, 21(2020), 2020.

[10] Rudy Bunel, Ilker Turkaslan, Philip HS Torr, Pushmeet Kohli, and M Pawan Kumar. A unified view of piecewise linear neural network verification. *Advances in Neural Information Processing Systems*, 2018.

[11] Alessandro De Palma, Harkirat Singh Behl, Rudy Bunel, Philip H. S. Torr, and M. Pawan Kumar. Scaling the convex barrier with active sets. In *International Conference on Learning Representations*, 2021.

[12] Alessandro De Palma, Harkirat Singh Behl, Rudy Bunel, Philip H. S. Torr, and M. Pawan Kumar. Scaling the convex barrier with sparse dual algorithms. *arXiv preprint arXiv:2101.05844*, 2021.

[13] Alessandro De Palma, Rudy Bunel, Alban Desmaison, Krishnamurthy Dvijotham, Pushmeet Kohli, Philip HS Torr, and M Pawan Kumar. Improved branch and bound for neural network verification via lagrangian decomposition. *arXiv preprint arXiv:2104.06718*, 2021.

[14] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting adversarial attacks with momentum. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9185–9193, 2018.

[15] Krishnamurthy Dvijotham, Robert Stanforth, Sven Gowal, Timothy Mann, and Pushmeet Kohli. A dual approach to scalable verification of deep networks. *Conference on Uncertainty in Artificial Intelligence*, 2018.

[16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[17] P. Henriksen and A. Lomuscio. Efficient neural network verification via adaptive refinement and adversarial search. In *Proceedings of the 24th European Conference on Artificial Intelligence (ECAI20)*, 2020.

[18] P. Henriksen and A. Lomuscio. Deepsplit: An efficient splitting method for neural network verification via indirect effect analysis. In *Proceedings of the 30th International Joint Conference on Artificial Intelligence (IJCAI21)*, To appear, August 2021.

[19] Florian Jaeckle and M Pawan Kumar. Generating adversarial examples with graph neural networks. *Conference on Uncertainty in Artificial Intelligence*, 2021.

[20] Florian Jaeckle, Jingyue Lu, and M Pawan Kumar. Neural network branch-and-bound for neural network verification. *arXiv preprint arXiv:2107.12855*, 2021.

[21] Guy Katz, Clark Barrett, David L Dill, Kyle Julian, and Mykel J Kochenderfer. Reluplex: An efficient smt solver for verifying deep neural networks. In *International Conference on Computer Aided Verification*, pages 97–117. Springer, 2017.

[22] Guy Katz, Derek A Huang, Duligur Ibeling, Kyle Julian, Christopher Lazarus, Rachel Lim, Parth Shah, Shantanu Thakoor, Haoze Wu, Aleksandar Zeljić, et al. The marabou framework for verification and analysis of deep neural networks. In *International Conference on Computer Aided Verification*, pages 443–452. Springer, 2019.

[23] P. Kouvaros and A. Lomuscio. Towards scalable complete verification of relu neural networks via dependency-based branching. In *Proceedings of the 30th International Joint Conference on Artificial Intelligence (IJCAI21)*, To Appear, 2021.

[24] Tim Kraska, Alex Beutel, Ed H Chi, Jeffrey Dean, and Neoklis Polyzotis. The case for learned index structures. In *Proceedings of the 2018 International Conference on Management of Data*, 2018.

[25] Changliu Liu, Tomer Arnon, Christopher Lazarus, and Mykel J Kochenderfer. Neuralverification.jl: Algorithms for verifying deep neural networks. In *ICLR 2019 Debugging Machine Learning Models Workshop*, 2019.

[26] Changliu Liu, Tomer Arnon, Christopher Lazarus, Christopher Strong, Clark Barrett, and Mykel J. Kochenderfer. Algorithms for verifying deep neural networks. *Foundations and Trends® in Optimization*, 4(3-4):244–404, 2021.

[27] Jingyue Lu and M Pawan Kumar. Neural network branching for neural network verification. In *International Conference on Learning Representations*, 2020.

[28] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *Proc. International Conference on Learning Representations (ICLR)*, 2018.

[29] Martin Vechev Mislav Balunovic. Adversarial training and provable defenses: Bridging the gap. In *Proc. International Conference on Learning Representations (ICLR)*, 2020.

[30] Mark Niklas Müller, Gleb Makarchuk, Gagandeep Singh, Markus Püschel, and Martin Vechev. Prima: Precise and general neural network certification via multi-neuron convex relaxations. *arXiv preprint arXiv:2103.03638*, 2021.

[31] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. *NIPS Autodiff Workshop*, 2017.

[32] Hadi Salman, Greg Yang, Huan Zhang, Cho-Jui Hsieh, and Pengchuan Zhang. A convex relaxation barrier to tight robustness verification of neural networks. *Advances in Neural Information Processing Systems*, 32:9835–9846, 2019.

[33] François Serre, Christoph Müller, Gagandeep Singh, Markus Püschel, and Martin Vechev. Scaling polyhedral neural network verification on GPUs. In *Proc. Machine Learning and Systems (MLSys)*, 2021.

[34] David Shriver, Sebastian G. Elbaum, and Matthew B. Dwyer. DNNV: A framework for deep neural network verification. In Alexandra Silva and K. Rustan M. Leino, editors, *Computer Aided Verification - 33rd International Conference, CAV 2021, Virtual Event, July 20-23, 2021, Proceedings, Part I*, volume 12759 of *Lecture Notes in Computer Science*, pages 137–150. Springer, 2021.

[35] David Shriver, Sebastian G. Elbaum, and Matthew B. Dwyer. Reducing DNN properties to enable falsification with adversarial attacks. In *43rd IEEE/ACM International Conference on Software Engineering, ICSE 2021, Madrid, Spain, 22-30 May 2021*, pages 275–287. IEEE, 2021.

[36] Gagandeep Singh, Rupanshu Ganvir, Markus Püschel, and Martin Vechev. Beyond the single neuron convex barrier for neural network certification. In *Advances in Neural Information Processing*

*Systems 32*, pages 15098–15109. Curran Associates, Inc., 2019.

[37] Gagandeep Singh, Timon Gehr, Matthew Mirman, Markus Püschel, and Martin Vechev. Fast and effective robustness certification. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 10802–10813. Curran Associates, Inc., 2018.

[38] Gagandeep Singh, Timon Gehr, Markus Püschel, and Martin Vechev. An abstract domain for certifying neural networks. *Proc. ACM Program. Lang.*, 3(POPL):41:1–41:30, 2019.

[39] Gagandeep Singh, Timon Gehr, Markus Püschel, and Martin Vechev. Boosting robustness certification of neural networks. In *International Conference on Learning Representations (ICLR)*. 2019.

[40] Gagandeep Singh, Markus Püschel, and Martin Vechev. Fast polyhedra abstract domain. In *Proc. Principles of Programming Languages (POPL)*, pages 46–59, 2017.

[41] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, pages 6105–6114. PMLR, 2019.

[42] Vincent Tjeng, Kai Y. Xiao, and Russ Tedrake. Evaluating robustness of neural networks with mixed integer programming. In *ICLR*, 2019.

[43] Hoang-Dung Tran, Stanley Bak, Weiming Xiang, and Taylor T. Johnson. Verification of deep convolutional neural networks using imagestars. In *32nd International Conference on Computer-Aided Verification (CAV)*. Springer, July 2020.

[44] Hoang-Dung Tran, Patrick Musau, Diego Manzanas Lopez, Xiaodong Yang, Luan Viet Nguyen, Weiming Xiang, and Taylor T. Johnson. Parallelizable reachability analysis algorithms for feed-forward neural networks. In *Proceedings of the 7th International Workshop on Formal Methods in Software Engineering (FormaliSE'19)*, FormaliSE '19, pages 31–40, Piscataway, NJ, USA, May 2019. IEEE Press.

[45] Hoang-Dung Tran, Patrick Musau, Diego Manzanas Lopez, Xiaodong Yang, Luan Viet Nguyen, Weiming Xiang, and Taylor T. Johnson. Star-based reachability analysis for deep neural networks. In *23rd International Symposium on Formal Methods (FM'19)*. Springer International Publishing, October 2019.

[46] Hoang-Dung Tran, Xiaodong Yang, Diego Manzanas Lopez, Patrick Musau, Luan Viet Nguyen, Weiming Xiang, Stanley Bak, and Taylor T. Johnson. NNV: The neural network verification tool for deep neural networks and learning-enabled cyber-physical systems. In *32nd International Conference on Computer-Aided Verification (CAV)*, July 2020.

[47] Joseph A. Vincent and Mac Schwager. Reachable polyhedral marching (rpm): A safety verification algorithm for robotic systems with deep neural network components, 2021.

[48] VNN-COMP. International verification of neural networks competition (VNN-COMP). *Verification of Neural Networks workshop at the International Conference on Computer-Aided Verification*, 2020.

[49] Shiqi Wang, Huan Zhang, Kaidi Xu, Xue Lin, Suman Jana, Cho-Jui Hsieh, and Zico Kolter. Beta-CROWN: Efficient bound propagation with per-neuron split constraints for complete and incomplete neural network verification. *arXiv preprint arXiv:2103.06624*, 2021.

[50] Eric Wong and Zico Kolter. Provable defenses against adversarial examples via the convex outer adversarial polytope. 2018.

[51] Haoze Wu, Alex Ozdemir, Aleksandar Zeljic, Kyle Julian, Ahmed Irfan, Divya Gopinath, Sadjad Fouladi, Guy Katz, Corina Pasareanu, and Clark Barrett. Parallelization techniques for verifying neural networks. In *# PLACEHOLDER_PARENT_METADATA_VALUE#*, volume 1, pages 128–137. TU Wien Academic Press, 2020.

[52] W. Xiang, H. Tran, and T. T. Johnson. Output reachable set estimation and verification for multi-layer neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 29(11):5777–5783, 2018.

[53] Kaidi Xu, Zhouxing Shi, Huan Zhang, Yihan Wang, Kai-Wei Chang, Minlie Huang, Bhavya

Kailkhura, Xue Lin, and Cho-Jui Hsieh. Automatic perturbation analysis for scalable certified robustness and beyond. *Advances in Neural Information Processing Systems*, 33, 2020.

[54] Kaidi Xu, Huan Zhang, Shiqi Wang, Yihan Wang, Suman Jana, Xue Lin, and Cho-Jui Hsieh. Fast and Complete: Enabling complete neural network verification with rapid and massively parallel incomplete verifiers. In *International Conference on Learning Representations*, 2021.

[55] Huan Zhang, Tsui-Wei Weng, Pin-Yu Chen, Cho-Jui Hsieh, and Luca Daniel. Efficient neural network robustness certification with general activation functions. *Advances in Neural Information Processing Systems*, 31:4939–4948, 2018.