

Systematic Generation of Diverse Benchmarks for DNN Verification

Dong Xu^[0000-0001-5643-7197], David Shriver^[0000-0003-0208-6517], Matthew Dwyer^[0000-0002-1937-1544], and Sebastian Elbaum^[0000-0001-9592-1352]

University of Virginia, Charlottesville VA 22904,
{dx3yy,dls2fc,matthewbdwyer,selbaum}@virginia.edu



Abstract. The field of verification has advanced due to the interplay of theoretical development and empirical evaluation. Benchmarks play an important role in this by supporting the assessment of the state-of-the-art and comparison of alternative verification approaches. Recent years have witnessed significant developments in the verification of deep neural networks, but diverse benchmarks representing the range of verification problems in this domain do not yet exist. This paper describes a neural network verification benchmark generator, GDVB, that systematically varies aspects of problems in the benchmark that influence verifier performance. Through a series of studies, we illustrate how GDVB can assist in advancing the sub-field of neural network verification by more efficiently providing richer and less biased sets of verification problems.

Keywords: Neural Network, Verification, Benchmark, Covering Array

1 Motivation

Advances in machine learning have enabled training of deep neural networks (DNN) that are capable of realizing complex functions that rival or exceed the performance of human-built software, e.g., [32,27,41]. This success has led system developers to deploy, or consider deployment of, DNN models in critical systems, e.g., [12,53,39]. Consequently, the verification of correctness properties of DNNs has become a key challenge to assuring autonomous systems, and the research community has risen to this challenge. In the three years since Katz et al. [30] presented RELUPLEX at CAV 2017, researchers have published more than 20 DNN verification approaches supporting different properties and DNN architectures and spanning a range of algorithmic approaches [20,30,63,22,46,50,56,9,19,62,45,59,60,61,29,13,18,14,36,31]. While DNN verification has its own unique challenges, it is also a recent example in the long-history of domain-specific verification research, e.g., for hardware[25], software[17], real-time systems[58], and cryptographic protocols[40], and can benefit from the experience of these communities.

A key lesson learned by the community is that despite the fact that verification emphasizes the development of theoretical and algorithmic techniques, *advances in verification research often arise from understanding how different algorithmic*

and implementation approaches compare – a process that requires empirical study. Empirical study in verification is common, but unlike many other fields of computer science, for decades it has organized *verification tool competitions* that serve as a regular and long-running form of community-driven empirical study. Researchers tracked the progress of SMT solvers over a span of 6 years at these community-driven empirical studies and found that repeatedly “a certain solver presents a key idea that improves the performance in a particular division, and this idea is implemented by most solvers” in the following year [7]. Enabling the type of comparative studies that drive such advances requires *verification benchmarks* – a fact that the verification community has recognized for at least 25 years, e.g., [55,33,43,10,8].

Benchmarking in verification has evolved in response to the demands of empirical study within the field, e.g., [3,4,1,2], to support two objectives: (A1) *assessment of the state-of-the-art* and (A2) *comparison of alternative approaches*. In support of these, the verification community has favored benchmarks that: (R1) **are diverse in structure and difficulty**; (R2) **represent verifier use cases**; and (R3) **evolve as verification technology advances**.

The verification benchmarking and competition literature suggests that these requirements are widely accepted. For example, the TPTP benchmark’s stated goals include R1 (“contains problems varying in difficulty”), R2 (“spans a diversity of subject matters”), and R3 (“is up-to-date”, “provides a mechanism for adding new problems”)[54]. Moreover, these requirements are promoted, either explicitly or implicitly, by many of the regularly held verification competitions. To meet R1 and R2 SAT competitions construct benchmarks that include problems from six different domains: software, hardware, A.I, obstruction, combinatorial challenges, and theorem proving[4]. SAT competitions since 2017 have instituted a *bring your own benchmarks* policy that requires verifier developers to submit 20 new benchmarks with at least 10 that are “not too easy” or “too hard” – which helps to address R1 and R3. SMT competitions have used selection criteria that are biased towards these same requirements, e.g., “balancing the difficulty of benchmarks”[7].

Verification competitions have undoubtedly been a positive force for developing high-quality verification benchmarks, but prior to their existence researchers were forced to develop their own “benchmarks” – a collection of verification problems on which they evaluate their techniques and perhaps others. This is the situation that the subfield of *DNN verification* finds itself in.

The risk in letting technique developers choose their own benchmark is selection bias – that the selected problems do not represent a broad or important population of problems. For example, if an SMT benchmark were selected based on the constraints generated by symbolic execution tools they would be structurally biased, consisting only of conjunctive formula. As another example, if a SAT benchmark were generated randomly it is likely that a large portion of the benchmark would not represent realistic use cases.

Good benchmarks are expensive to develop, e.g., [11], but they are an invaluable resource for advancing a research community. When well designed they seek to balance requirements R1-R3 and to support a fair and accurate assessment of

the state-of-the-art and comparison between alternative algorithmic and implementation approaches. This paper reports on GDVB, the first *framework for systematic Generation of DNN Verification problem Benchmarks*, that meets the de-facto requirements for verification benchmarks, R1-R3, in order to support objectives A1-A2 for the rapidly evolving field of DNN verification.

GDVB takes a **generative** approach to benchmark development – an approach that has risen in popularity in recent years [35,5,64]. Unlike, other generative benchmark approaches GDVB seeks to systematically cover variations in verification problems that are known to influence verifier performance. Towards that end, GDVB is parameterized by: (1) a set of *factors* known to influence the performance of DNN verifiers; (2) a *coverage* goal that determines the combination of factors that should be reflected in the benchmark; and (3) a *seed* verification problem from which a set of variant problems are generated. From these parameters, it computes a constrained mixed-level covering array[15] defining a set of factor-value tuples. Each tuple defines how the seed verification problem can be transformed to give rise to a verification problem capable of exposing performance variation in a DNN verifier.

As a benchmark generator GDVB naturally meets requirement R3. By starting from a seed network representing a DNN verification use case, GDVB is guaranteed to meet R2. As we discuss in §4, the use of factors allows GDVB to produce systematically diverse verification problems both in terms of structure and difficulty in order to meet requirement R1. Moreover, GDVB offers the potential to reduce selection bias in performing evaluations of DNN verifiers, since it assures coverage of a space of performance related factors. Finally, GDVB is designed to support the rapidly evolving field of DNN verifiers by allowing the generation of benchmarks, e.g., from new seeds as verifiers improve, as new performance factors are identified, and to target challenge problems in different DNN domains, e.g., regression models for autonomous UAV navigation [39,53].

The contributions of this paper are: identification of the need for unbiased and diverse benchmarks for DNN verification; a study of factors that affect the performance of DNN verification tools (§3); the specification of a verification benchmark as the solution to a constrained mixed-level covering array problem (§4); the GDVB algorithm for computing a benchmark from a verification problem by transforming the neural network and correctness specification (§4.3); the evaluation of GDVB on multiple state-of-the-art DNN verifiers using different seed verification problems that demonstrates how GDVB results can support the evaluation of DNN verifiers (§5); and the GDVB tool.

2 Background and Related Work

Deep Neural Networks (DNN) A DNN is trained to accurately approximate a target function, $f : \mathbb{R}^d \rightarrow \mathbb{R}^r$. A network, $n : \mathbb{R}^d \rightarrow \mathbb{R}^r$, is comprised of a graph of L hidden layers, l_1, \dots, l_L , along with an input layer, $l_{in} = l_0$, and output layer, $l_{out} = l_{L+1}$. Each hidden layer defines an independent function, where their

composition when applied to the output of l_{in} generates values in l_{out} that define the network output.

Hidden layers are, generally, comprised of a set of *neurons* that accumulate a weighted sum of their inputs from the prior layer and then apply an *activation function* to determine how to non-linearly scale that sum to compute the output from the layer. A variety of different activation functions have been explored in the literature, including: rectified linear units (ReLU), sigmoid, and tanh.

The design of a DNN involves choosing an appropriate set of *layer types*, e.g., convolutional, maxpooling, fully-connected, the instantiation of those layers, e.g., the number of neurons, the specific activation function, and the definition of how layers are interconnected. Together these comprise the DNN *architecture* [23].

Networks are trained using a variety of algorithmic strategies with the goal of minimizing the loss in the approximation of the learned function relative to some proxy for f , e.g., labeled training data. The training process is stochastic, e.g., initial weight values are randomized, which leads to variation in n even when architecture, training algorithm, and training data are fixed.

§3 reveals how DNN architecture can influence verification performance.

DNN Specifications Given a network $n : \mathbb{R}^d \rightarrow \mathbb{R}^r$, a property, ϕ , defines a set of constraints over the inputs, ϕ_x , and an associated set of constraints over the outputs, ϕ_y . Verification of n seeks to prove: $\forall \mathbf{x} \in \mathbb{R}^d : \phi_x(\mathbf{x}) \Rightarrow \phi_y(\mathbf{N}(x))$ where $\mathbf{N}(x)$ is running the neural network n with input x .

Specifying behavioral properties of DNNs is challenging and is an active area of research [24]. In [30], a set of 188 purely conjunctive properties, of the form described above, were defined for a simple neural network, with 7 inputs, encoding of a rule set for autonomous aircraft collision avoidance (ACAS). In [44,60,59], properties expressing output range invariants were used, for example, that the steering angle never exceeded an absolute value of 30 degrees. Much of the work on DNN verification has focused on local robustness properties [51,50,52], which state that for a selected target input the output of the network is invariant for other inputs within a specified distance of the target.

§3 reveals how the specification can influence verification performance.

DNN Verification Methods and Tools There are a variety of different algorithmic and implementation approaches taken to verifying the validity of a DNN with respect to a stated correctness property.

Definition 1. A DNN verification problem, $\langle n, \phi \rangle$, is comprised of a DNN, n , and a property specification, ϕ . The outcome of a verification problem for a DNN verifier indicates whether $n \models \phi$ is valid, invalid, or unknown – indicating that the problem cannot be determined to be either valid or invalid.

A recent DNN verification survey [37], classifies approaches as being based on reachability, optimization, and search algorithms – or their combination. Reachability methods begin with a symbolic encoding of an input set and compute, for each layer, a symbolic encoding of the output set. They vary in the symbolic encodings used, e.g., intervals, polyhedra, and in the degree of overapproximation they introduce [63,22,46,50]. Optimization methods formulate verification as an

optimization problem whose solution implies the validity of ϕ [38,56,9,19,62,45]. Search methods work in combination with reachability and optimization, by decomposing the input space to formulate verification sub-problems that are discharged by the above techniques [60,61,29,13,30,14,59,18,20].

In this paper, we use implementations of the following verifiers: ERAN [50], BAB [14], NEURIFY [59], PLANET [20], and RELUPLEX [30].

Verification Benchmarking We covered the broad landscape of work on benchmark development for verification in (§1). There have been efforts to develop benchmarks within a variety of different verification problem domains, e.g. hardware[25], software[17], real-time systems[58], cryptographic protocols[40], and for different encodings of verification problems, e.g., model checking [33], SAT [4], SMT [8], and theorem proving [55].

In recent work on DNN verification, researchers have shared collections of examples that, in a sense, serve as informal benchmarks and permit comparative evaluation, e.g. [30,50]. While valuable, these examples were not intended to, and do not, comprise a benchmark meeting requirements R1-R3. To our knowledge, GDVB is the first approach to achieving those goals for DNN verification.

For several years, the SAT community has been exploring scalable benchmarks, e.g., [35,21]. For instance, to explore conflict-driven clause learning (CDCL) SAT solver performance, Elffers et al. [21] used crafted parameterized benchmarks that can be scaled with respect to different factors that may influence performance. We conduct a similar domain analysis of factors, but focus on the landscape of DNN verification algorithms developed to date. Like this line of work, GDVB advocates a scalable approach to benchmark generation. As described in §4, GDVB starts with seed problems that are challenging for current verifiers and “scales them down”, but it can also be applied to start with easier seed problems and “scale them up” as more typical of the prior work on scalable benchmarking.

Verification Benchmark Ranking The verification community has explored a variety of ranking schemes for assessing the cost-effectiveness of techniques. A key challenge is that verification techniques vary not only in their cost, e.g., time to produce a verification result, but also in their accuracy, e.g., whether they produce an *unknown* result. For example, SAT competitions have employed a range of scoring models, e.g., purse-based ranking, *solution-count ranking* (SCR), careful ranking, and penalized average runtime (PAR2) [6]. SCR, which counts the number of solved problem instances and uses verification time as a tie breaker [57], is the scoring system of choice [4,1]. In §5, we report DNN verifier performance using both SCR and PAR2 scoring systems.

Covering Arrays In §3 we explore factors that influence DNN verifier performance. Studying all their combinations would be cost prohibitive, so we consider weaker notions of coverage.

A covering array defines a systematic method for testing how combinations of parameter values influence system performance [16]. A covering array is an $N \times k$ array. The k columns represent *factors* that may influence performance and cells can take on v *levels* – defining settings for factors. The N rows of the array define combinations of factor-levels. Arrays are defined to achieve a *strength*

of the coverage, t . $t = 2$ defines pairwise strength, which means that all pairs of levels for all factors are present in some row of the covering array.

We require a richer form of covering array that permits the number of levels to vary with different factors, i.e., a mixed-level covering array (MCA), and that can constrain specified factor-level combinations, e.g., by forbidding their inclusion in the MCA. By modeling each factor as a variable and its levels as the domain of the variable, one can express constraints as propositional logic formulae over equality terms; if the levels are ordered then richer underlying theories can be applied. A constrained-MCA defines an MCA that is consistent with a given constraint, C .

Definition 2. *Constrained Mixed-level Covering Array (Def. 2.9 from [15])*
 $CMCA(N; t, k, (|v_1|, |v_2|, \dots, |v_k|), C)$ is an $N \times k$ array on $|v|$ symbols, where $|v| = \sum_{i=0}^k |v_i|$, with the following properties: 1) Each column i ($1 \leq i \leq k$) contains only elements from a set S_i of size $|v_i|$, 2) the rows of each $N \times t$ subarray cover all t -tuples of values from the t columns at least one time, and 3) all rows are models of C .

Transforming Neural Networks The GDVB approach manipulates factors that influence DNN verifier performance to construct a diverse benchmark. For DNN construction, we leverage a recent approach, R4V [47], that given an original DNN and an architectural specification automates the transformation of the DNN and uses distillation[28] to train it to closely match the test accuracy of the original DNN. R4V transformation specifications can be written to change a number of architectural parameters of a network including: the input dimension, the range of values for each input dimension, the number of layers, the number of neurons per layer, the number of convolutional kernels, and the stride and padding of a convolutional layer.

3 Identifying Factors that Influence Verifier Performance

As discussed in §1 the verification community has acted to create policies that incentivize *diverse* benchmarks. Diversity is desirable in a benchmark because it (a) demonstrates the range of applicability of a verification technology and (b) exposes performance variation within and across verification technologies. Consider, that the SMT competition benchmark selection process seeks to “include equal numbers of satisfiable and unsatisfiable benchmarks at different levels of difficulty”[7]. This is due to the fact that the SMT community understands that the satisfiability or unsatisfiability of a benchmark problem is a factor that influences verifier performance¹.

GDVB seeks to make factors influencing verifier performance explicit and to manipulate them to generate a diverse benchmark. To determine an initial set of factors for DNN verifiers we began with an analysis of the literature, which

¹ Since unsatisfiability requires the consideration of all possible variable assignments which generally is more costly than finding a single satisfiable assignment.

identified several candidate factors, and then conducted a targeted and exploratory **factor study** to identify whether *manipulating a factor could influence some performance measure of some DNN verifier*. This study only aims to identify such factors and does not seek to characterize the complex relationship between factors and DNN verifier performance; for example, we do not aim to capture a comprehensive set of factors, assess the independence of or relations between factors, or rank factors in terms of their degree of influence. A richer and more detailed factor study might further improve the utility of GDVB, but we leave such a study to future work.

3.1 Potential Factors

Relatively few published papers on DNN verification explicitly discuss the factors that influence performance, but nearly all of them present metrics on the verification problems they solved.

Evaluation results for RELUPLEX present data on verifier outcome and solve time for local robustness properties that vary in the input center point and radius [30]; most subsequent papers report similar property variation. Evaluation results for ROBUSTVERIFIER present a study of varying the number of layers in the DNN and its impact on verifier performance[36]. Evaluation results for ERAN present performance variations across a range of networks varying in the number of layers, layer types, and neurons[22,51,50,52]. Bunel et al. [14] were the first that we are aware of to explicitly vary factors of DNN verification problems. They found that the performance varied with input dimension, number of neurons per layer, and number of layers across a set of 6 different DNN verifiers. All of the other papers published on DNN verification in recent years have used verification problems that varied, in an ad-hoc fashion, over a subset of the above factors.

3.2 Exploratory Factor Study

As in other verification domains, DNN verifier performance is multi-faceted. In our study, we consider both verification time and accuracy. We say that the result of a verification problem is *accurate* if a verifier determines conclusively that the problem is *valid* or *invalid*, result as opposed to *unknown*².

We study factors associated with both properties and DNNs. Based on the literature analysis, we identified 2 factors related to the correctness property: *scale* and *translation*. Scaling a property involves increasing the size of the input domain which will involve *more DNN behavior* in verification. Translating a property involves moving it to a different location in the input domain which will involve *different DNN behavior* in verification. For robustness properties, scaling and translation involve changing the radius and center point of the hypercube describing the input space under verification. One might wonder whether rotation of a property can influence verification performance. For robustness properties,

² We cross-check accurate results with multiple verifiers.

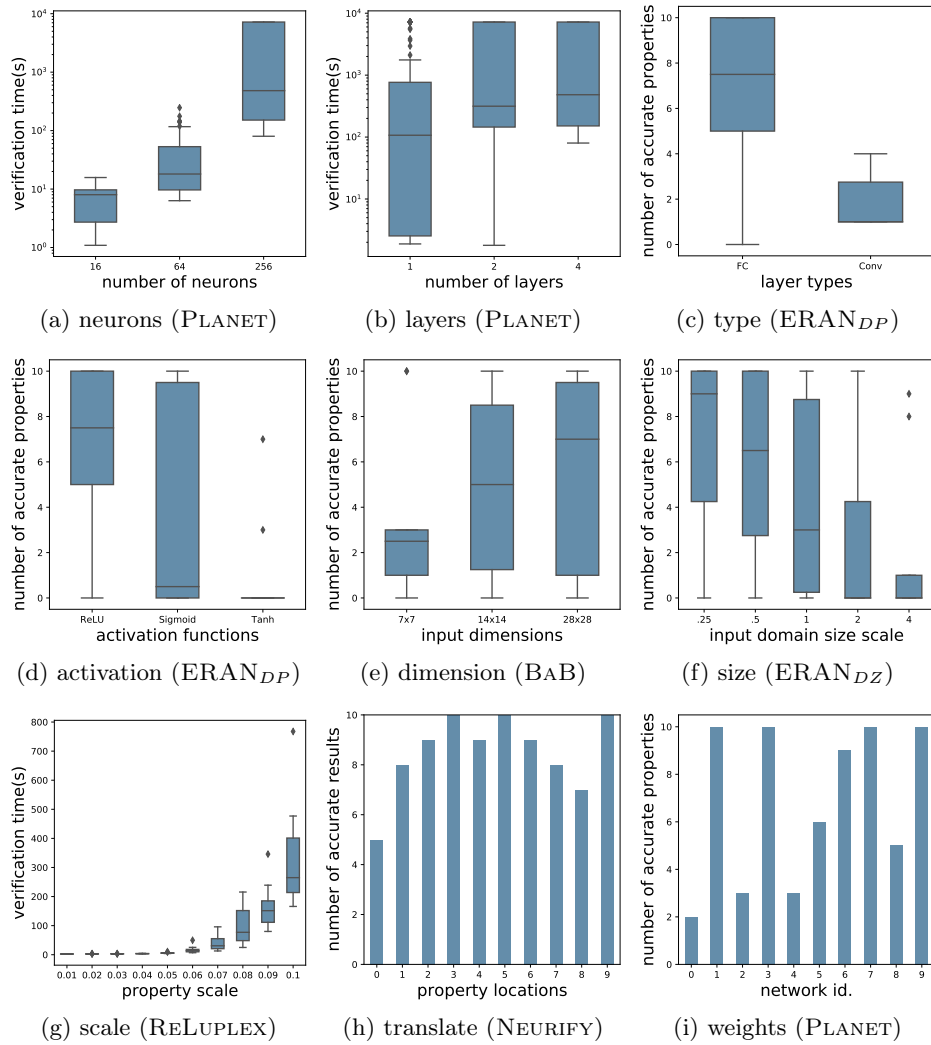


Fig. 1. DNN Verifier Performance Across Factors

this seems unlikely given their symmetry, but it could be a factor for more irregular input regions – we leave this for future work.

Based on the literature analysis, we identified 4 factors related to the DNN: number of *neurons*, number of *layers*, the *type* of layers, the input *dimension*. We conjectured that an additional 3 factors might impact verifier performance: the type of *activation* function, the input domain *size*, and the learned *weights*.

Our exploratory factor study is opportunistic in that we seek to find a verification problem for which manipulation of a selected factor exhibits performance variation. Towards this end, we conducted a series of trials where we vary a factor hypothesized to influence verification performance, while holding all other factors constant and report the results in Fig. 1. We studied variations of networks for

the MNIST task and considered local robustness properties since these were well-supported across a range of different verifiers. We used different verifiers across the study: RELUPLEX, PLANET, NEURIFY, BAB, ERAN with the DeepPoly (DP) and DeepZono (DZ) abstract domains. We now briefly describe the trials and then summarize the outcome.

Number of Neurons: The architecture of the DNN was fixed, with 4 fully-connected layers using ReLU activation functions, and the total number of neurons was varied (16, 64, 256) – they were spread evenly across layers. Each network is trained 10 times and verified on 100 local robustness properties. Fig. 1(a) plots the number of neurons versus verification time for PLANET. **Verification time can increase with the number of neurons.**

Number of Layers: We use the same context as for the neuron factor study, except that we fixed the number of neurons at 256 and vary the number of layers (1,2,4). Fig. 1(b) plots the number of layers versus verification time for PLANET. **Verification time can increase with the number of layers .**

Layer Types: We use a pair of two-layer neural networks, with the same number of neurons, where one has a fully-connected layer and the other a convolutional layer. Each network is trained 10 times and verified on 10 local robustness properties. Fig. 1(c) plots layer type versus the number of properties for which accurate results are produced using ERAN_{DP}. **Verification accuracy can vary with layer type.**

Activation Function: We use the fully-connected network from the layer types study, we generated three networks by altering the activation function to use sigmoid and tanh. The training setup and properties remain the same as in the previous trial. Fig. 1(d) plots the activation function versus the number of properties for which accurate results are produced using ERAN_{DP}. **Verification accuracy can vary with activation function.**

Input Dimension: We use 3 architectures that differ only in their input dimension which is scaled ($\frac{1}{16}, \frac{1}{4}, 1$) relative on the original problem. The training setup and properties are from the layer type study. Fig. 1(e) plots the input dimension versus the number of properties for which accurate results are produced using BAB. **Verification accuracy can increase with increasing input dimension.**

Input Size: We use 5 architectures that differ only in the range of values of their inputs which are scaled ($\frac{1}{4}, \frac{1}{2}, 1, 2, 4$) based on the original problem. The training setup and properties are from the layer type study. Fig. 1(f) plots the input size versus the number of properties for which accurate results are produced using ERAN_{DZ}. **Verification accuracy can decrease with increasing input domain size.**

Property Scale: We use a single-layer network and reuse the training setup and properties from the layer type study. We scale the properties (0.01 – 0.1) to generate verification problems. Fig. 1(g) plots property scaling versus the verification time using RELUPLEX. **Verification time can increase with increasing property scale.**

Property Translation: We replicated the property scale study, but held the scale fixed and translated the center point of the local robustness property to 10 other locations. Fig. 1(h) plots the number of DNNs for each of the 10 translated properties for which accurate results could be produced using NEURIFY. **Verification accuracy can vary with property translation.**

Network Weights: Building of the property studies, we explore the verification of 10 scaled property variants across the same network trained 10 times with different initial weights. Fig. 1(i) plots the number of accurate properties for which the results could be produced using PLANET. **Verification accuracy can vary with the learned weights of the network.**

Exploratory Study Findings Varying the factors studied influences the performance of different DNN verifiers differently – in terms of time or accuracy. For example, we found that: varying input dimension impacts BAB’s accuracy, but not RELUPLEX’s; varying input domain size impacts ERAN_{DZ}’s accuracy, but not NEURIFY’s; and varying property scale impacts RELUPLEX’s verification time, but not NEURIFY’s.

This study provides a starting set of viable factors that can be used to parameterize the GDVB approach to produce verification problem benchmarks in which those factors are systematically varied. Furthermore, as we discuss in §4, GDVB generative process allows for us to accommodate information about new factors that might be revealed in future factor studies.

4 The GDVB Approach

The goal of GDVB is to meet requirements R1-R3 by producing a *factor diverse* benchmark that (a) reflects aspects of the complexity encoded in a real verification problem that acts as a seed for generation $\langle n_s, \phi_s \rangle$, (b) varies aspects of the problem that are related to verifier performance, (c) accounts for interactions among those factors, and (d) is comprised only of well-defined verification problems.

Rather than synthesize random verification problems, we seed the generation process in order to generate a benchmark that reflects the complexity of the seed problem. This permits benchmarks to be generated to reflect the challenges present in different DNN problem sub-domains.

Factors, like those described in §3, may interact; changes to one factor may mask or amplify DNN verifier performance changes arising from another. Exploring all combinations of factors is expensive, but by using covering arrays we can systematically explore interactions among factors. Accounting for such interactions helps to produce a benchmark that is *less biased* than one that only covers individual factor variations.

Not all combinations of factors are possible. For example, if one reduces the number of layers in a network to 0, then it is not possible to preserve the number of neurons in the original network. Thus, benchmark generation must take into account constraints among factors to ensure that only well-defined problems are included in a benchmark.

4.1 Factor Diverse Benchmarks

Consider a set of factors, F , with a set of levels, L_f , for each factor, $f \in F$; we refer to L_f as the *level set* of f . For a verification problem, p , let $l(p)$ be the set of factor levels corresponding to the problem. A benchmark, B , is a set of verification problems and we can denote the factor levels for the benchmark as $l(B) = \{l(p) \mid p \in B\}$.

The simplest form of diversity for a benchmark is requiring that all individual factor levels be present in at least one verification problem, $\forall f \in F : \forall l \in L_f : \exists p \in l(B) : l \in p$. However, this diversity fails to account for interactions among factors. The simplest form of interaction-sensitive diversity considers pairs of factors, but as we discuss below our approach generalizes to any arity of factor-level coverage.

For a pair of factors, $f, f' \in F$, the Cartesian product of their level sets defines the set of all pairwise combinations of their levels. Across all factors the set of such pairs is $\text{pairs}(F) = \{(l, l') \mid f, f' \in F \wedge f \neq f' \wedge l \in L_f \wedge l' \in L_{f'}\}$. A *pairwise diverse benchmark* is one in which

$$\forall (x, y) \in \text{pairs}(F) : \exists p \in l(B) : (x, y) \in \{(x', y') \mid x' \in p \wedge y' \in p\}$$

Constraints on allowable combinations of factors serve to restrict a benchmark. A pairwise exclusion constraint, $\gamma(F) \subseteq \text{pairs}(F)$, requires that

$$\forall (x, y) \in \gamma(F) : \forall p \in l(B) : \neg(x \in p \wedge y \in p)$$

We write γ when F is understood from the context.

The arity of factor-level coverage and exclusion constraints can vary independently. It is common for factor-level coverage to be uniform and to generalize it to t -way coverage, i.e., to require coverage of the elements of the Cartesian product of the level sets of t factors. On the other hand, as observed in prior work [15], constraints generally involve a mix of arity. To denote this generality we define $\Gamma \subseteq \bigcup_i \gamma_i$ where γ_i defines the set of possible i -way exclusion constraints.

Example. Consider the DAVE-2 DNN which accepts 100 by 100 color images and infers an output indicating the steering angle[12]. DAVE-2 is comprised of 5 convolutional layers with 55296, 17424, 3888, 3136, and 1600 neurons, respectively, followed by 4 fully connected layers with 1164, 100, 50, and 10 neurons, respectively. All 82668 neurons use ReLU activations. One can define a local robustness property for DAVE-2 as

$$\phi = \forall \mathbf{x} \in i \pm 0.02 : \|\text{DAVE-2}(\mathbf{x}) - \text{DAVE-2}(i)\| \leq 5$$

which states that for a given an input image, i , all inputs within a distance of 0.02 will result in an inferred steering angle within 5 degrees of the angle for i . These yield the verification problem $\langle \text{DAVE-2}, \phi \rangle$.

Consider factors for the number of neurons, number of convolutional layers, and number of fully-connected layers; a tuple $(\#neuron, \#conv, \#fc)$ represents levels for these factors. For each factor consider two percentage levels: 100%

and 50%. A neuron factor level of 50% indicates that a version of DAVE-2 with 41334 neurons is required. In the absence of constraints, an example pairwise factor diverse benchmark for $\langle \text{DAVE-2}, \phi \rangle$ consists of the following four verification problems: (100%, 100%, 100%), (100%, 50%, 50%), (50%, 100%, 50%), and (50%, 50%, 100%). The property ϕ is constant across the benchmark.

4.2 From Factor Covering Arrays to Verification Problems

Given a set of factors, $F = \{f_1, f_2, \dots, f_{|F|}\}$, and levels, L_{f_i} , a t -way factor diverse benchmark of k verification problems is specified by

$$\text{CMCA}(|F|; t, k, (|L_{f_1}|, |L_{f_2}|, \dots, |L_{f_{|F|}}|), \Gamma)$$

Each element in this mixed level covering array specifies how to construct a verification problem in the benchmark from the seed problem.

Levels are operationalized as transformations on verification problems. We assume a sufficient set of transformations, Δ , such that a verification problem can be transformed into a form that achieves any level of any factor

$$\forall f \in F : \forall l_f \in L_f : \exists \delta \in \Delta : l_f \in l(\delta(\langle n_s, \phi_s \rangle))$$

The definition of Δ and L_i must be coordinated to achieve this property.

A per-factor transformation $\delta \in \Delta$ may impact a single component of a verification problem, e.g., reducing the number of neurons in a DNN does not impact the property, or both components, e.g., the input dimension impacts the DNN and the property by transforming the input data domain. The set of all transformations Δ defines the set of verification problems that can be produced by application of a set of per-factor transformations to the seed problem,

$$\Delta(\langle n_s, \phi_s \rangle) = \{ \langle n, \phi \rangle \mid \langle n, \phi \rangle = \delta_{f_1} \circ \delta_{f_2} \dots \circ \delta_{f_{|F|}}(\langle n_s, \phi_s \rangle) \wedge \delta_i \in \Delta \}$$

The set of all possible factor level combinations is $\prod_{f \in F} L_f$, i.e., the product of all of the per-factor levels. The set of t -way factor level combinations is

$$c_t = \{ c \mid a \in \prod_{f \in F} L_f \wedge c \subseteq a \wedge |c| = t \}$$

allowing for the interpretation of $|F|$ -tuples as sets.

Definition 3. *Given a set of factors F , with associated factor levels L_f , a t -way factor diverse benchmark, B , for a seed problem $\langle n_s, \phi_s \rangle$ with exclusion constraints Γ is defined by the following: (1) $B \subseteq \Delta(\langle n_s, \phi_s \rangle)$; (2) $\forall \langle n, \phi \rangle \in B : \forall \gamma \in \Gamma : \gamma \not\subseteq l(\langle n, \phi \rangle)$; and (3) $\forall c \in c_t - \Gamma : \exists \langle n, \phi \rangle \in B : c \subseteq l(\langle n, \phi \rangle)$*

4.3 Generating Benchmarks

GDVB is defined in Algorithm 1. We use existing techniques, e.g. Automated Combinatorial Testing for Software (ACTS) [34], for generating a CMCA for constraints specified as logical formulae where factors are variables and levels are

Algorithm 1: GDVB($\langle n_s, \phi_s \rangle, F, \Gamma, t$) Algorithm

Data: a seed problem $\langle n_s, \phi_s \rangle$, a set of factors F and constraints Γ , a coverage goal t
Result: A benchmark of DNN verification problems B

- 1 $C \leftarrow \text{GENCMCA}(F, \Gamma, t)$
- 2 $B \leftarrow \emptyset$
- 3 **for** $c \in C$ **do**
- 4 $B \leftarrow B \cup \text{TRANSFORM}(\langle n_s, \phi_s \rangle, c)$
- 5 **end**

values for those variables. A CMCA is a set of k -tuples. Each such tuple defines the target level for each factor for a problem in the generated benchmark. Those levels are used to transform the given seed verification problem and the resultant problem is accumulated in the benchmark.

TRANSFORM uses different approaches to transform the seed DNN and the property. DNN transformation builds on an approach called R4V that automates architectural transformations to DNNs by scaling (1) the number of neurons in a fully connected layer, (2) the number of kernels in a convolutional layer, (3) the input dimension, or (4) the range of values within an input dimension [47]. The first 3 of these require changes to the structure of the DNN and the last two require changes to the training data, e.g., reshaping, renormalizing. R4V ensures that the network is well-defined after transformation. TRANSFORM maps factor-levels to per-layer scale parameters for R4V.

R4V permits the training of a network using network distillation which we find advantageous for GDVB because: it accelerates the training process, and it drives training to match the accuracy of the problem DNN to that of n_s , which reduces variation in accuracy across B . We adapt R4V so that after each training epoch, the learned DNN weights and the validation accuracy is recorded. When training finishes, we select the weights associated with the highest validation accuracy. Training is performed using the training data and hyperparameters for n_s .

Whereas R4V can be used to directly manipulate DNN architecture related factors, it can only indirectly affect the learned weights. To address this, we adopt the approach taken throughout the machine learning literature – train a network on multiple initial seeds and report performance across seeds. Thus, each DNN in B is trained multiple times, thereby producing a benchmark comprised of $s * |B|$ verification problems, where s is the desired number of seeds.

DNN Transformation Example. Consider this element of the CMCA described above: $\langle (50\%, 100\%, 50\%), \phi \rangle$, applied to DAVE-2. TRANSFORM would compute that 50% of the fully connected layers should be present in the resultant DNN and randomly select 2 of the 4 layers to scale by 0. The fully-connected layers are chosen at random, since the layer count factor does not consider layer ordering. If we consider the case where the layers with 100 and 50 neurons are dropped, this will eliminate 150 neurons. The other transformation required is to reduce the number of neurons by half. To do that all remaining layers will be scaled by $\frac{82668 * 0.5 - 150}{82668} = 0.498$.

Property transformation builds on a domain-specific language (DSL) for specifying DNN correctness properties defined by the *deep neural network verification framework* (DNNV) [48]. Specifications in this Python-based DSL are parametric and TRANSFORM maps factor-levels to those parameters. For example, Fig. 2 defines the parametric local robustness property ϕ that is centered at the image stored at “path/to/image”, has radius 0.02, and can be translated and scaled through parameters \mathbf{t} and \mathbf{s} , respectively.

Restricting factors to levels that are supported by TRANSFORM and using CMCA algorithms that meet Def. 2 ensures that GDVB produces a solution that meets Def. 3.

```

N=Network("N")
s=Parameter("s",float,
            default=1.0)
e=0.02*s
x=Image("path/to/image")
t=np.load(Parameter("t",
                    str,
                    "path/to/zeros.npy"))
x=x+t
forall(x_,
    Implies(
        (x-e)<x_<(x+e),
        abs(N(x_)-N(x)) <= 5
    )
)

```

Fig. 2. Parametric Property ϕ

4.4 An Instantiation of GDVB

We developed an instance of GDVB³ that supports a set of factors informed by the results of the study in §3, *percentage-based levels* for those factors, and a set of constraints that restrict benchmark problems to those that are non-trivial and that can be efficiently trained

Our instantiation of GDVB supports the following factors: the total number of neurons in the DNN (**neu**), the number of fully-connected layers (**fc**), the number of convolutional layers (**conv**), the dimension of the DNN input (**idm**), the size of each DNN input dimension (**ids**), the scale of the property (**scl**), and the translation of the property (**trn**). We do not support an activation function factor because only ERAN support non-ReLU activations and, thus, using them would render other verifiers inapplicable for large portions generated benchmarks.

We use quintile factor levels, {20%, 40%, 60%, 80%, 100%}, for factors neu, idm, ids, and scl. To permit the elimination of layer types we extend these levels with an additional quintile, 0%, for fc and conv. For trn, we select a set of five translations that shift the property to be centered on a different instance of the training data; unlike the above levels this level is unordered.

Our instantiation of GDVB exclusion constraints for DAVE-2 are as follows: (1) $fc = 0 \wedge conv = 0$, (2) $conv = 0 \wedge neu \geq 20$, (3) $conv = 0 \wedge idm \geq 80$, and (4) $conv = 100 \wedge idm = 20$. The first of these requires that some layer be present. The second and third are related to the blowup in the size of fully-connected layers that results from dropping all convolutional layers which makes training difficult; limiting the total number of neurons and the reduction input dimension mitigates this. The fourth constraint ensures that the input dimension reduction results in a meaningful network; without it the dimensionality reduction achieved by sequences of convolutional layers yields an invalid network, i.e., the input to some layer is smaller than the kernel size.

³ <https://github.com/edwardxu0/GDVB>

These constraints were developed iteratively based on feedback from the R4V tool, which reports when TRANSFORM has specified an invalid DNN, and when training failed to closely approximate the accuracy of the seed network.

We note that this instance of GDVB is flexible in that it permits the customization of levels, as we demonstrate in the next section, to generate a benchmark that focuses on variation in a subset of factors. More generally, GDVB can easily be extended to support additional factors and levels for which an instance of TRANSFORM can be defined. We expect that GDVB will evolve in this way as studies of DNN verifiers are performed.

5 GDVB in Use

In this section we showcase the potential uses of GDVB across a series of artifacts and verifiers, while highlighting the challenges it helps to systematically address.

5.1 Setup

Our evaluation applies GDVB to two seed networks: $MNIST_{ConvBig}$ and DAVE-2. We selected $MNIST_{ConvBig}$ because it is one of the largest networks in ERAN’s evaluation [50]; it includes 4 convolutional layers and 3 fully connected layers with 48,074 neurons and 1,974,762 parameters. We selected DAVE-2 to illustrate the application of GDVB to a larger network that has been the subject of other DNN analysis [42]; it has 5 convolutional layers and 5 fully connected layers with 82,669 neurons and 2,116,983 parameters.

Table 1 lists the 9 verifiers we selected for our study. This list includes the most well-known verifiers and verification algorithms. We also select variations of some verification approaches. We use Branch-and-Bound (BAB), as well as a variation of Branch-and-Bound with Smart-Branching (BABSB). Additionally, we evaluate the ERAN verifier with 4 available abstract domains: DeepZono ($ERAN_{DZ}$), DeepPoly ($ERAN_{DP}$), RefineZono ($ERAN_{RZ}$), and RefinePoly ($ERAN_{RP}$).

Table 1. Verifiers used in GDVB study

Verifier	Algorithm
RELUPLEX [30]	Search-Optimization
PLANET [20]	Search-Optimization
BAB [14]	Search-Optimization
BABSB [14]	Search-Optimization
NEURIFY ⁴ [59]	Optimization
$ERAN_{DZ}$ [50]	Reachability
$ERAN_{DP}$ [51]	Reachability
$ERAN_{RZ}$ [52]	Reachability
$ERAN_{RP}$ [49]	Reachability

To evaluate verifier performance, we use the *solution-count ranking* (SCR)[57], which counts the number of properties that returned accurate verification results. Additionally, we measured the *penalized average runtime* (PAR2)[6], which is computed as the sum of the verification times for *sat* and *unsat* results and twice time limit for all other verification results.

⁴ We use the version of NEURIFY provided in DNNV[48], which is modified to be applicable to a wide range of problems, whereas the original version was hard-coded to a particular verification problem[59].

Verifier	MNIST _{ConvBig}		DAVE-2	
	SCR	PAR2	SCR	PAR2
ERAN _{DZ}	11.40±0.49	18,126.80±488.27	7.20±1.94	24,496.20±1,176.59
ERAN _{DP}	21.00±0.89	9,206.00±806.70	18.40±2.15	17,443.00±1,344.65
ERAN _{RZ}	10.20±0.40	19,252.60±343.66	5.80±2.14	25,236.60±1,253.90
ERAN _{RP}	12.60±1.02	16,981.40±930.71	10.20±1.83	22,250.60±1,186.44
NEURIFY	22.00±1.10	8,636.20±1,008.63	19.20±2.56	17,247.80±1,397.05
PLANET	7.00±0.63	23,145.60±468.18	3.40±1.62	27,268.60±775.56
BAB	0.20±0.40	28,689.80±220.40	0.00±0.00	28,800.00±0.00
BASB	0.00±0.00	28,800.00±0.00	0.00±0.00	28,800.00±0.00
RELUPLEX	3.20±0.40	25,757.80±381.40	4.40±1.02	26,023.60±635.90

Table 2. Mean & Variance of SCR and PAR2 Scores Across Benchmarks. (The darker and lighter gray boxes indicate the best and second best results.)

All training and verification took place under CentOS Linux 7. R4V transformation and distillation jobs ran on NVIDIA 1080Ti GPUs. Verification jobs were limited to 4 hours and ran on 2.3GHz and 2.2GHz Xeon processors with 64GB of memory, for DAVE-2 and MNIST_{ConvBig}, respectively.

5.2 Comparing verifiers across a range of challenges

Consider the use case where a researcher is attempting to compare a new verifier (e.g., a new algorithm, a revised implementation, an extension to an existing approach) against existing verifiers. As shown earlier, for such comparison to be meaningful, many factors must be considered and properly explored. Given a seed network, a property, a set of factors, and a coverage goal, GDVB can generate a benchmark that helps to reduce bias in conducting such an evaluation.

For this use case we consider seed networks and local robustness properties similar to those from the ERAN_{DZ} study [50] for the MNIST_{ConvBig} verification problem and local robustness properties based on those from the NEURIFY study [59] for the DAVE-2 verification problem. We run an instance of GDVB using the factors and levels described in §4.4, a coverage strength of 2, and train 5 versions of each network to account for stochastic weight variation. The total time to generate and train GDVB (MNIST_{ConvBig}, ...) was 24.3 hours and the resulting 30 verification problems took 401.8 hours to run across all 9 verifiers. For GDVB (DAVE-2, ...) 44 verification problems were generated with training and verification times of 158.2 hours and 772.4 hours, respectively. CMCA generation took less than a minute for both problems. Each problem in the benchmark must be trained and verified in sequence, but across problems they can be parallelized. We exploited this to reduce the cost of running the benchmarks to 4.9 hours for MNIST_{ConvBig} and 7.9 hours for DAVE-2. We measured the SCR and PAR2 score for the nine verifiers across the benchmarks.

The results are shown in Table 2. Since the SCR and PAR2 score trends are the same we depict just SCR in Fig. 3. Boxplots show the SCR scores for a verifier across all the generated problems; variation in plots arises from the 5

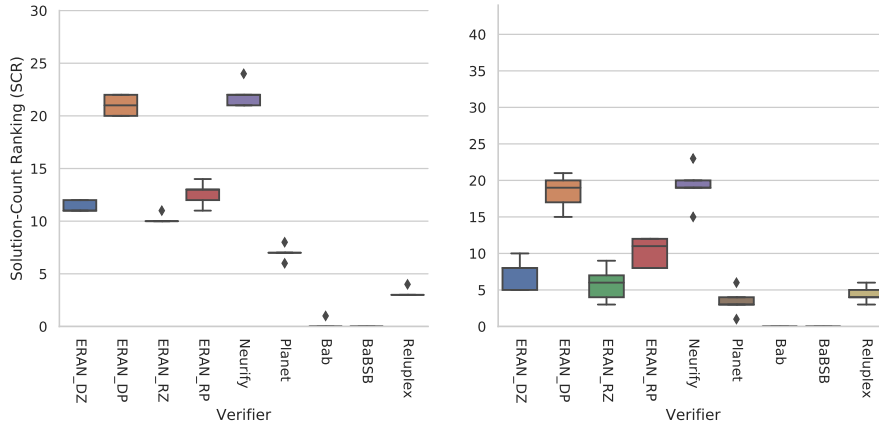


Fig. 3. SCR Score for Nine Verifiers on GDVB Benchmarks with $\text{MNIST}_{ConvBig}$ (left) and DAVE-2 (right) seeds

trained versions of the networks for each problem. For each box, the middle line represent the median, the box-bounds are the first and third quartiles, and the whiskers represent minimal and maximal values.

The plot for $\text{MNIST}_{ConvBig}$ on the left of Fig. 3 shows that **the GDVB benchmark with the $\text{MNIST}_{ConvBig}$ seed is able to identify considerable performance variation across verifiers**, with ERAN_{DP} and NEURIFY accurately verifying a median of over 20 properties, the rest of the ERAN -variants verifying between 10 and 13 properties, and the remaining tools verifying between 0 and 8 properties. The results are consistent when we employ DAVE-2 as the seed network, with **marked differences among groups of verifiers** although the generated problems turned out to be more challenging across all verifiers. ERAN_{DP} and NEURIFY , the top performers, can verify less than half of the generated problems. Verifiers like BAB were unable to verify any problem derived from DAVE-2 because of the complexity of the seed problem. This point highlights the need for benchmarks to evolve with networks that incorporate emerging technology, and also GDVB’s ability to automatically generate a benchmark from different seeds to address that need.

Now, understanding the overall performance of a family of verifiers is useful but it is likely just the first step for a researcher to understand under what conditions a verifier excels or struggles. When such conditions correspond to the factors manipulated by GDVB, then they are readily available for further analysis. One analysis may consist of simply plotting the data across its multiple dimensions. We do so in the form of radar-charts for DAVE-2 in Fig. 4 and for $\text{MNIST}_{ConvBig}$ in Fig. 5⁵. Since the observations we can gather from both networks are similar, we just discuss DAVE-2 in detail. Each chart includes six axes representing a factor scaled between 0 and 1. The solid lines link the

⁵ We do not plot BABS8 as its performance was identical to BAB .

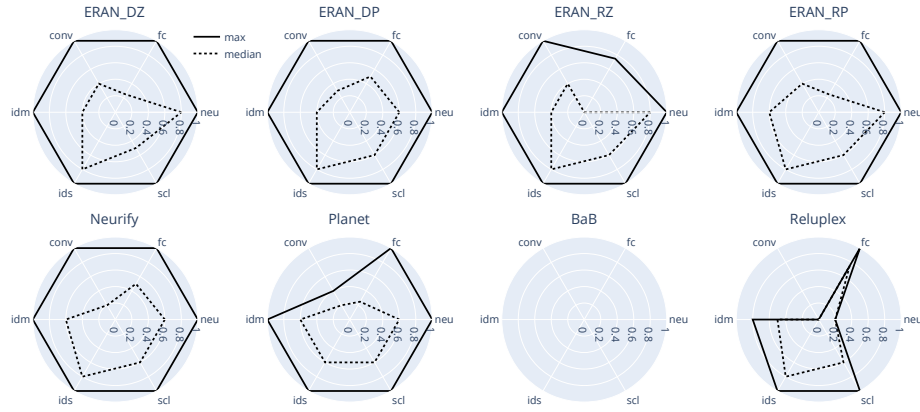


Fig. 4. DAVE-2: Radar plot with maximum (solid) and median (dotted) values

maximum values across factors that were accurately verified while the dotted lines link the median values across factors.

The shape of the lines in the radar plots clearly show that the **verification problems generated by GDVB reveal unique patterns across the verifiers**. For example, the RELUPLEX plot indicates that it can do well verifying networks with multiple fully connected (fc) layers but is challenged by larger networks (neu) and those with convolutional layers (conv). Comparing multiple charts also reveals some interesting trade-offs. For example, for smaller networks with just fully connected layers, the medians seem to indicate that RELUPLEX is better than PLANET. However, when a network incorporates convolutional layers or a larger number of neurons, PLANET appears to outperform RELUPLEX.

Looking across charts can also pinpoint specific improvements resulting from tool extensions or revisions. For example, the median line of $ERAN_{RZ}$ indicates that it was not as effective in handling verification problems with a larger number of layers as its predecessor $ERAN_{DZ}$; the same trend holds for the pair $ERAN_{RP}$ and $ERAN_{DP}$. We note that a more restrictive benchmark that is biased towards fewer fully connected layers might not reveal such differences.

GDVB offers the opportunity to investigate such differences even further by generating targeted verification problems for a subset of factors hypothesized to be culprits of those differences. For example, GDVB could generate additional verification problems with a number of fully connected layers between 60% and 80% of the total, while keeping the other factors constant, to refine the understanding of the differences between $ERAN_{RZ}$ and $ERAN_{DZ}$.

This study illustrates how GDVB benchmarks support the exploration of verifier performance, lowering the burden on researchers to manually prepare tens to hundreds of verification problems, and reducing the opportunities for bias.

5.3 GDVB and benchmark requirements R1-R3

As explained in §1, benchmarking in verification seeks to develop benchmarks that are: diverse; representative of real use cases; and reactive to new technologies.

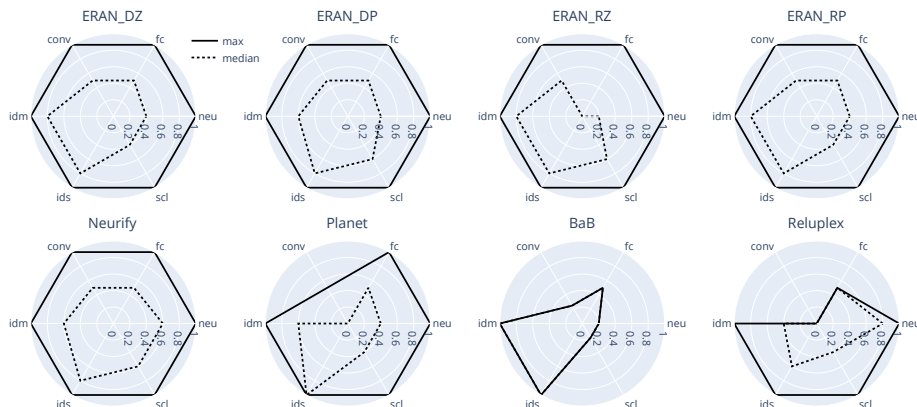


Fig. 5. $MNIST_{ConvBig}$: Radar plot with maximum (solid) and median (dotted) values

The previous sections have provided evidence of how, through its generative nature, GDVB is reactive to new advances in technology included in the seed network. We have also seen the high degree of parameterization GDVB offers including for setting a seed network from which realistic attributes are inherited in the generated verification problems. In this section we want to illustrate how GDVB addresses the diversity requirement.

To depict diversity we use the parallel coordinate graph in Fig. 6. Each vertical line corresponds to a factor, and the markers in each vertical line corresponds to an explored level. Each verification problem is a polyline that connects the factors’ levels explored by it. The two sets of lines correspond to the verification problems included in the DAVE-2 benchmark published with NEURIFY [59], which is a downsized version of the full DAVE-2 DNN, and the benchmark produced by GDVB (DAVE-2, ...). Each factor in the plot is normalized by dividing by the maximum value for the factor.

Fig. 6 shows that the NEURIFY’s DAVE-2 has a large number of neurons, inputs, and dimensions. Yet, it provides very limited coverage of all the factor levels that may affect verification performance. In contrast, GDVB provides a systematic exploration of the factors levels that can affect verifier performance making it much less biased – especially to the numbers of layers in the verification problems, and the combination of those factor levels.

The parallel plot for GDVB benchmark with the $MNIST_{ConvBig}$ seed (not shown for space reasons), depicts a similar trend in terms of systematic exploration of diversity, but since $MNIST_{ConvBig}$ is simpler than DAVE-2, the generated benchmark is correspondingly simpler. This points to the need to identify representative and challenging seeds when parameterizing GDVB. GDVB is fully capable of accomodating factor levels that exceed 100% of a seed network, which is a means of pushing verifiers to the limits of their abilities.

We note that excluding factors or levels can yield a systematically generated benchmark that is unable to characterize differences between verifiers, or worse, misleads such a characterization by emphasizing certain factors while overlooking others. For example, not exploring different network sizes or exploring networks

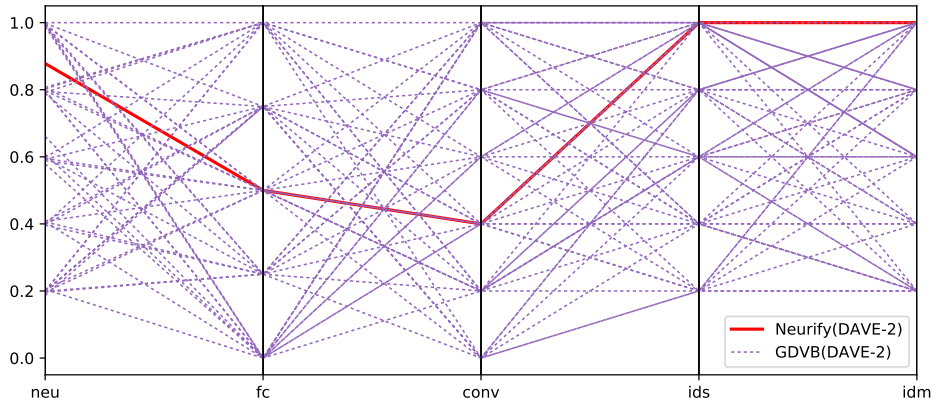


Fig. 6. Diversity explored across factor levels

sizes under 1000 neurons will render similar scores across many DNN verifiers that are differentiated by more comprehensive benchmarks. In applying GDVB, we suggest selecting as many factors as we know may matter, starting from a challenging seed problem, and incrementally refining the levels as needed to focus benchmark results to differentiate verifier performance.

6 Conclusion

The increasing adoption of DNNs has led to a surge in research on DNN verification techniques. Benchmarks to assess these emerging techniques, however, are costly to develop, often lack in diversity and do not represent the population of real evolving DNNs. To address this challenge, we have introduced GDVB, a framework for systematically generating DNN verification problems seeded in complex, real-world networks, ensuring that benchmarks are derived from real problems. GDVB is parameterizable by the factors that may influence verification performance and thereby supports scalable benchmarking. A preliminary study, using 9 DNN verifiers, demonstrates how GDVB can support the assessment of the state-of-the-art.

We plan to conduct broader studies of verifier performance using GDVB, and we encourage other researchers to use and contribute to it. There are many directions to explore in identifying new factors that influence performance, e.g., the impact of quantization and model compression approaches [26]. Work in this direction promises to deepen the community’s understanding and lead to advances in DNN verification.

Acknowledgements

This material is based in part upon work supported by the National Science Foundation under grant numbers 1901769 and 1900676, by the U.S. Army Research Office under grant number W911NF-19-1-0054.

References

1. 14th International Satisfiability Modulo Theories Competition, <https://smt-comp.github.io/2019/>
2. Competition on Software Verification, <https://sv-comp.sosy-lab.org/2019/>
3. Hardware Model Checking Competition, <http://fmv.jku.at/hwccc19/index.html>
4. The International Satisfiability Competitions, <http://www.satcompetition.org/>
5. Amendola, G., Ricca, F., Truszczynski, M.: A generator of hard 2QBF formulas and ASP programs. In: 16th International Conference on Principles of Knowledge Representation and Reasoning (2018)
6. Balint, A., Belov, A., Järvisalo, M., Sinz, C.: Overview and analysis of the SAT challenge 2012 solver competition. *Artificial Intelligence* **223**, 120–155 (2015)
7. Barrett, C., Deters, M., De Moura, L., Oliveras, A., Stump, A.: 6 years of SMT-COMP. *Journal of Automated Reasoning* **50**(3), 243–277 (2013)
8. Barrett, C., Stump, A., Tinelli, C.: The SMT-LIB standard: Version 2.0 **13**, 14 (2010)
9. Bastani, O., Ioannou, Y., Lampropoulos, L., Vytiniotis, D., Nori, A.V., Criminisi, A.: Measuring neural net robustness with constraints. In: Proceedings of the 30th International Conference on Neural Information Processing Systems. pp. 2621–2629 (2016)
10. Beyer, D., Löwe, S., Wendler, P.: Reliable benchmarking: requirements and solutions. *International Journal on Software Tools for Technology Transfer* **21**(1), 1–29 (2019)
11. Blackburn, S.M., Garner, R., Hoffmann, C., Khang, A.M., McKinley, K.S., Bentzur, R., Diwan, A., Feinberg, D., Frampton, D., Guyer, S.Z., Hirzel, M., Hosking, A., Jump, M., Lee, H., Moss, J.E.B., Phansalkar, A., Stefanović, D., VanDrunen, T., von Dincklage, D., Wiedermann, B.: The DaCapo benchmarks: Java benchmarking development and analysis. In: Proceedings of the 21st annual ACM SIGPLAN conference on Object-oriented programming systems, languages, and applications. pp. 169–190 (2006)
12. Bojarski, M., Testa, D.D., Dworakowski, D., Firner, B., Flepp, B., Goyal, P., Jackel, L.D., Monfort, M., Muller, U., Zhang, J., Zhang, X., Zhao, J., Zieba, K.: End to end learning for self-driving cars. *CoRR* (2016)
13. Boopathy, A., Weng, T.W., Chen, P.Y., Liu, S., Daniel, L.: CNN-Cert: An efficient framework for certifying robustness of convolutional neural networks. In: Association for the Advancement of Artificial Intelligence (Jan 2019)
14. Bunel, R., Turkaslan, I., Torr, P.H., Kohli, P., Kumar, M.P.: A unified view of piecewise linear neural network verification. In: Proceedings of the 32nd International Conference on Neural Information Processing Systems. pp. 4795–4804 (2018)
15. Cohen, M.B., Dwyer, M.B., Shi, J.: Constructing interaction test suites for highly-configurable systems in the presence of constraints: A greedy approach. *IEEE Transactions on Software Engineering* **34**(5), 633–650 (Sep 2008)
16. Cohen, M.B., Gibbons, P.B., Mugridge, W.B., Colbourn, C.J.: Constructing test suites for interaction testing. In: 25th International Conference on Software Engineering. pp. 38–48 (May 2003)
17. D’silva, V., Kroening, D., Weissenbacher, G.: A survey of automated techniques for formal software verification. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* **27**(7), 1165–1178 (2008)
18. Dutta, S., Jha, S., Sankaranarayanan, S., Tiwari, A.: Output range analysis for deep feedforward neural networks. In: NASA Formal Methods Symposium. pp. 121–138 (2018)

19. Dvijotham, K., Stanforth, R., Goyal, S., Mann, T., Kohli, P.: A dual approach to scalable verification of deep networks. In: Proceedings of the 34th Conference Annual Conference on Uncertainty in Artificial Intelligence. pp. 162–171 (2018)
20. Ehlers, R.: Formal verification of piece-wise linear feed-forward neural networks. In: Automated Technology for Verification and Analysis. pp. 269–286 (2017)
21. Elffers, J., Giráldez-Cru, J., Gocht, S., Nordström, J., Simon, L.: Seeking practical CDCL insights from theoretical SAT benchmarks. In: International Joint Conferences on Artificial Intelligence. pp. 1300–1308 (2018)
22. Gehr, T., Mirman, M., Drachler-Cohen, D., Tsankov, P., Chaudhuri, S., Vechev, M.: AI2: Safety and robustness certification of neural networks with abstract interpretation. In: IEEE Symposium on Security and Privacy. pp. 3–18 (May 2018)
23. Goodfellow, I., Bengio, Y., Courville, A., Bengio, Y.: Deep learning, vol. 1 (2016)
24. Gopinath, D., Converse, H., Pasareanu, C.S., Taly, A.: Property inference for deep neural networks. In: 34th IEEE/ACM International Conference on Automated Software Engineering. pp. 797–809 (2019)
25. Gupta, A.: Formal hardware verification methods: A survey. In: Computer Aided Verification. pp. 5–92 (1992)
26. Han, S., Mao, H., Dally, W.J.: Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding (2015)
27. Hinton, G., Deng, L., Yu, D., Dahl, G.E., Mohamed, A.r., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T.N.: Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. IEEE Signal processing magazine **29**(6), 82–97 (2012)
28. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. In: NIPS Deep Learning and Representation Learning Workshop (2015)
29. Huang, X., Kwiatkowska, M., Wang, S., Wu, M.: Safety verification of deep neural networks. In: Computer Aided Verification - 29th International Conference. pp. 3–29 (2017)
30. Katz, G., Barrett, C.W., Dill, D.L., Julian, K., Kochenderfer, M.J.: Reluplex: An efficient SMT solver for verifying deep neural networks. In: Computer Aided Verification - 29th International Conference. pp. 97–117 (2017)
31. Katz, G., Huang, D.A., Ibeling, D., Julian, K., Lazarus, C., Lim, R., Shah, P., Thakoor, S., Wu, H., Zeljić, A., et al.: The Marabou framework for verification and analysis of deep neural networks. In: International Conference on Computer Aided Verification. pp. 443–452 (2019)
32. Krizhevsky, A., Sutskever, I., Hinton, G.E.: ImageNet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems 25: Proceedings of the 26th Annual Conference on Neural Information Processing Systems 2012. pp. 1106–1114 (2012)
33. Kropf, T.: Benchmark-circuits for hardware-verification. In: International Conference on Theorem Provers in Circuit Design. pp. 1–12 (1994)
34. Kuhn, R., Kacker, R.: Automated Combinatorial Testing for Software, <https://csrc.nist.gov/projects/automated-combinatorial-testing-for-software>
35. Lauria, M., Elffers, J., Nordström, J., Vinyals, M.: CNFgen: A generator of crafted benchmarks. In: Theory and Applications of Satisfiability Testing. pp. 464–473 (2017)
36. Lin, W., Yang, Z., Chen, X., Zhao, Q., Li, X., Liu, Z., He, J.: Robustness verification of classification deep neural networks via linear programming. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 11418–11427 (2019)

37. Liu, C., Arnon, T., Lazarus, C., Barrett, C., Kochenderfer, M.J.: Algorithms for verifying deep neural networks. *CoRR* (2019)
38. Lomuscio, A., Maganti, L.: An approach to reachability analysis for feed-forward ReLU neural networks. *CoRR* (2017)
39. Loquercio, A., Maqueda, A.I., Blanco, C.R.D., Scaramuzza, D.: Dronet: Learning to fly by driving. *IEEE Robotics and Automation Letters* (2018)
40. Meadows, C.A.: Formal verification of cryptographic protocols: A survey. In: *International Conference on the Theory and Application of Cryptology*. pp. 133–150 (1994)
41. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., Hassabis, D.: Human-level control through deep reinforcement learning. *Nature* **518**(7540), 529–533 (Feb 2015)
42. Pei, K., Cao, Y., Yang, J., Jana, S.: DeepXplore: Automated whitebox testing of deep learning systems. In: *Proceedings of the 26th Symposium on Operating Systems Principles*. pp. 1–18 (2017)
43. Pelánek, R.: BEEM: benchmarks for explicit model checkers. In: *International SPIN Workshop on Model Checking of Software*. pp. 263–267 (2007)
44. Pulina, L., Tacchella, A.: An abstraction-refinement approach to verification of artificial neural networks. In: *International Conference on Computer Aided Verification*. pp. 243–257 (2010)
45. Raghunathan, A., Steinhardt, J., Liang, P.: Certified defenses against adversarial examples. In: *The International Conference on Learning Representations* (2018)
46. Ruan, W., Huang, X., Kwiatkowska, M.: Reachability analysis of deep neural networks with provable guarantees. In: *International Joint Conferences on Artificial Intelligence*. pp. 2651–2659 (2018)
47. Shriver, D., Xu, D., Elbaum, S.G., Dwyer, M.B.: Refactoring neural networks for verification. *CoRR* (2019)
48. Shriver, D.L.: Deep Neural Network Verification Toolbox, <https://github.com/dlshriver/DNNV>
49. Singh, G., Ganvir, R., Püschel, M., Vechev, M.: Beyond the single neuron convex barrier for neural network certification. In: Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., Garnett, R. (eds.) *Advances in Neural Information Processing Systems* 32, pp. 15072–15083 (2019)
50. Singh, G., Gehr, T., Mirman, M., Püschel, M., Vechev, M.: Fast and effective robustness certification. In: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (eds.) *Advances in Neural Information Processing Systems* 31, pp. 10802–10813 (2018)
51. Singh, G., Gehr, T., Püschel, M., Vechev, M.: An abstract domain for certifying neural networks. *Proceedings of the ACM on Programming Language* **3** (Jan 2019)
52. Singh, G., Gehr, T., Püschel, M., Vechev, M.: Boosting robustness certification of neural networks. In: *Proceedings of the International Conference on Learning Representations* (2019)
53. Smolyanskiy, N., Kamenev, A., Smith, J., Birchfield, S.: Toward low-flying autonomous MAV trail navigation using deep neural networks for environmental awareness. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems*. pp. 4241–4247 (Sep 2017)
54. Sutcliffe, G.: The TPTP problem library and associated infrastructure. *Journal of Automated Reasoning* **43**(4), 337–362 (2009)

55. Sutcliffe, G., Suttner, C.: The TPTP problem library. *Journal of Automated Reasoning* **21**(2), 177–203 (1998)
56. Tjeng, V., Xiao, K.Y., Tedrake, R.: Evaluating robustness of neural networks with mixed integer programming. In: *International Conference on Learning Representations* (2019)
57. Van Gelder, A.: Careful ranking of multiple solvers with timeouts and ties. In: Sakallah, K.A., Simon, L. (eds.) *Theory and Applications of Satisfiability Testing*. pp. 317–328 (2011)
58. Wang, F.: Formal verification of timed systems: A survey and perspective. *Proceedings of the IEEE* **92**(8), 1283–1305 (2004)
59. Wang, S., Pei, K., Whitehouse, J., Yang, J., Jana, S.: Efficient formal safety analysis of neural networks. In: *Advances in Neural Information Processing Systems*. pp. 6367–6377 (2018)
60. Wang, S., Pei, K., Whitehouse, J., Yang, J., Jana, S.: Formal security analysis of neural networks using symbolic intervals. In: *USENIX Security Symposium*. pp. 1599–1614 (2018)
61. Weng, T., Zhang, H., Chen, H., Song, Z., Hsieh, C., Daniel, L., Boning, D.S., Dhillon, I.S.: Towards fast computation of certified robustness for ReLU networks. In: *International Conference on Machine Learning. Proceedings of Machine Learning Research*, vol. 80, pp. 5273–5282 (2018)
62. Wong, E., Kolter, J.Z.: Provable defenses against adversarial examples via the convex outer adversarial polytope. In: *International Conference on Machine Learning. Proceedings of Machine Learning Research*, vol. 80, pp. 5283–5292 (2018)
63. Xiang, W., Tran, H., Johnson, T.T.: Output reachable set estimation and verification for multilayer neural networks. *IEEE Transactions on Neural Networks and Learning Systems* **29**(11), 5777–5783 (Nov 2018)
64. You, J., Wu, H., Barrett, C., Ramanujan, R., Leskovec, J.: G2SAT: Learning to generate SAT formulas. In: *Advances in Neural Information Processing Systems*. pp. 10552–10563 (2019)