

Exploiting Verified Neural Networks via Floating Point Numerical Error

Kai Jia¹ and Martin Rinard¹

MIT CSAIL, Cambridge MA 02139, USA
{jiakai,rinard}@mit.edu

Abstract. Researchers have developed neural network verification algorithms motivated by the need to characterize the robustness of deep neural networks. The verifiers aspire to answer whether a neural network guarantees certain properties with respect to all inputs in a space. However, many verifiers inaccurately model floating point arithmetic but do not thoroughly discuss the consequences.

We show that the negligence of floating point error leads to unsound verification that can be systematically exploited in practice. For a pre-trained neural network, we present a method that efficiently searches inputs as witnesses for the incorrectness of robustness claims made by a complete verifier. We also present a method to construct neural network architectures and weights that induce wrong results of an incomplete verifier. Our results highlight that, to achieve practically reliable verification of neural networks, any verification system must accurately (or conservatively) model the effects of any floating point computations in the network inference or verification system.

Keywords: Verification of neural networks · Floating point soundness · Tradeoffs in verifiers

1 Introduction

Deep neural networks (DNNs) have been successful at various tasks, including image processing, language understanding, and robotic control [30]. However, they are vulnerable to adversarial inputs [40], which are input pairs indistinguishable to human perception that cause a DNN to give substantially different predictions. This situation has motivated the development of network verification algorithms that claim to prove the robustness of a network [3, 33, 42], specifically that the network produces identical classifications for all inputs in a perturbation space around a given input.

Verification algorithms typically reason about the behavior of the network assuming real-valued arithmetic. In practice, however, the computation of both the verifier and the neural network is performed on physical computers that use floating point numbers and floating point arithmetic to approximate the underlying real-valued computations. This use of floating point introduces numerical error that can potentially invalidate the guarantees that the verifiers claim to

provide. Moreover, the existence of multiple software and hardware systems for DNN inference further complicates the situation because different implementations exhibit different numerical error characteristics. Unfortunately, prior neural network verification research rarely discusses floating point (un)soundness issues (Section 2).

This work considers two scenarios for a decision-making system relying on verified properties of certain neural networks: (i) The adversary can present arbitrary network inputs to the system while the network has been pretrained and fixed; (ii) The adversary can present arbitrary inputs and also network weights and architectures to the system. We present an efficient search technique to find witnesses of the unsoundness of complete verifiers under the first scenario. The second scenario enables inducing wrong results more easily, as will be shown in Section 5. Note that even though allowing arbitrary network architectures and weights is a stronger adversary, it is still practical. For example, one may deploy a verifier to decide whether to accept an untrusted network based on its verified robustness, and an attacker might manipulate the network so that its nonrobust behavior does not get noticed by the verifier.

Specifically, we train robust networks on the MNIST and CIFAR10 datasets. We work with the `MIPVerify` complete verifier [42] and several inference implementations included in the PyTorch framework [29]. For each implementation, we construct image pairs $(\mathbf{x}_0, \mathbf{x}_{\text{adv}})$ where \mathbf{x}_0 is a brightness-modified natural image, such that the implementation classifies \mathbf{x}_{adv} differently from \mathbf{x}_0 , \mathbf{x}_{adv} falls in a ℓ_∞ -bounded perturbation space around \mathbf{x}_0 , and the verifier incorrectly claims that no such adversarial image \mathbf{x}_{adv} exists for \mathbf{x}_0 within the perturbation space. Moreover, we show that if modifying network architecture or weights is allowed, floating point error of an incomplete verifier `CROWN` [49] can also be exploited to induce wrong results. Our method of constructing adversarial images is not limited to our setting but is applicable to other verifiers that do not soundly model floating point arithmetic.

We emphasize that any verifier that does not correctly or conservatively model floating point arithmetic fails to provide any safety guarantee against malicious network inputs and/or network architectures and weights. Ad hoc patches or parameter tuning can not fix this problem. Instead, verification techniques should strive to provide soundness guarantees by correctly incorporating floating point details in both the verifier and the deployed neural network inference implementation. Another solution is to work with quantized neural networks that eliminate floating point issues [20].

2 Background and related work

Training robust networks: Researchers have developed various techniques to train robust networks [24, 26, 43, 47]. Madry et al. [24] formulates the robust training problem as minimizing the worst loss within the input perturbation and proposes training on data generated by the Projected Gradient Descent

(PGD) adversary. In this work, we consider robust networks trained with the PGD adversary.

Complete verification: Complete verification (a.k.a. exact verification) methods either prove the property being verified or provide a counterexample to disprove it. Complete verifiers have formulated the verification problem as a Satisfiability Modulo Theories (SMT) problem [3, 12, 17, 21, 34] or a Mixed Integer Linear Programming (MILP) problem [6, 10, 13, 23, 42]. In principle, SMT solvers are able to model exact floating point arithmetic [32] or exact real arithmetic [8]. However, for efficiency reasons, deployed SMT solvers for verifying neural networks all use inexact floating point arithmetic to reason about the neural network inference. MILP solvers typically work directly with floating point, do not attempt to model real arithmetic exactly, and therefore suffer from numerical error. There have also been efforts on extending MILP solvers to produce exact or conservative results [28, 39], but they exhibit limited performance and have not been applied to neural network verification.

Incomplete verification: On the spectrum of the tradeoff between completeness and scalability, incomplete methods (a.k.a. certification methods) aspire to deliver more scalable verification by adopting over-approximation while admitting the inability to either prove or disprove the properties in certain cases. There is a large body of related research [11, 14, 26, 31, 37, 45, 46, 49]. Salman et al. [33] unifies most of the relaxation methods under a common convex relaxation framework and suggests that there is an inherent barrier to tight verification via layer-wise convex relaxation captured by such a framework. We highlight that floating point error of implementations that use a direct dot product formulation has been accounted for in some certification frameworks [36, 37] by maintaining upper and lower rounding bounds for sound floating point arithmetic [25]. Such frameworks should be extensible to model numerical error in more sophisticated implementations like the Winograd convolution [22], but the effectiveness of this extension remains to be studied. However, most of the certification algorithms have not considered floating point error and may be vulnerable to attacks that exploit this deficiency.

Floating point arithmetic: Floating point is widely adopted as an approximate representation of real numbers in digital computers. After each calculation, the result is rounded to the nearest representable value, which induces roundoff error. A large corpus of methods have been developed for floating point analysis [2, 9, 38, 41], but they have yet not been applied to problems at the scale of neural network inference or verification involving millions of operations. Concerns for floating point error in neural network verifiers are well grounded. For example, the verifiers **Reluplex** [21] and **MIPVerify** [42] have been observed to occasionally produce incorrect results on large scale benchmarks [15, 44]. However, no prior work tries to systematically invalidate neural network verification results via exploiting floating point error.

3 Problem definition

We consider 2D image classification problems. Let $\mathbf{y} = \text{NN}(\mathbf{x}; \mathbf{W})$ denote the classification confidence given by a neural network with weight parameters \mathbf{W} for an input \mathbf{x} , where $\mathbf{x} \in \mathbb{R}_{[0,1]}^{m \times n \times c}$ is an image with m rows and n columns of pixels each containing c color channels represented by floating point values in the range $[0, 1]$, and $\mathbf{y} \in \mathbb{R}^k$ is a logits vector containing the classification scores for each of the k classes. The class with the highest score is the classification result of the neural network.

For a logits vector \mathbf{y} and a target class number t , we define the Carlini-Wagner (CW) loss [5] as the score of the target class subtracted by the maximal score of the other classes:

$$L_{\text{CW}}(\mathbf{y}, t) = y_t - \max_{i \neq t} y_i \quad (1)$$

Note that \mathbf{x} is classified as an instance of class t if and only if $L_{\text{CW}}(\text{NN}(\mathbf{x}; \mathbf{W}), t) > 0$, assuming no equal scores of two classes.

Adversarial robustness of a neural network is defined for an input \mathbf{x}_0 and a perturbation bound ϵ , such that the classification result is stable within allowed perturbations:

$$\begin{aligned} \forall \mathbf{x} \in \text{Adv}_\epsilon(\mathbf{x}_0) : L(\mathbf{x}) > 0 \\ \text{where } L(\mathbf{x}) = L_{\text{CW}}(\text{NN}(\mathbf{x}; \mathbf{W}), t_0) \\ t_0 = \text{argmax NN}(\mathbf{x}_0; \mathbf{W}) \end{aligned} \quad (2)$$

In this work we consider ℓ_∞ -norm bounded perturbations:

$$\text{Adv}_\epsilon(\mathbf{x}_0) = \{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}_0\|_\infty \leq \epsilon \wedge \min \mathbf{x} \geq 0 \wedge \max \mathbf{x} \leq 1\} \quad (3)$$

We use the MIPVerify [42] complete verifier to demonstrate our attack method. MIPVerify formulates (2) as an MILP instance $L^* = \min_{\mathbf{x} \in \text{Adv}_\epsilon(\mathbf{x}_0)} L(\mathbf{x})$ that is solved by the commercial solver Gurobi [16]. The network is robust if $L^* > 0$. Otherwise, the minimizer \mathbf{x}^* encodes an adversarial image.

Due to the inevitable presence of numerical error in both the network inference system and the verifier, the exact specification of $\text{NN}(\cdot; \mathbf{W})$ (i.e., a bit-level accurate description of the underlying computation) is not clearly defined in (2). We consider the following implementations included in the PyTorch framework to serve as our candidate definitions of the convolutional layers in $\text{NN}(\cdot; \mathbf{W})$, while nonconvolutional layers use the default PyTorch implementation:

- $\text{NN}_{\text{C,M}}(\cdot; \mathbf{W})$: A matrix-multiplication-based implementation on x86/64 CPUs. The convolution kernel is copied into a matrix that describes the dot product to be applied on the flattened input for each output value.
- $\text{NN}_{\text{C,C}}(\cdot; \mathbf{W})$: The default convolution implementation on x86/64 CPUs.
- $\text{NN}_{\text{G,M}}(\cdot; \mathbf{W})$: A matrix-multiplication-based implementation on NVIDIA GPUs.
- $\text{NN}_{\text{G,C}}(\cdot; \mathbf{W})$: A convolution implementation using the IMPLICIT_GEMM algorithm from the cuDNN library [7] on NVIDIA GPUs.

- $\text{NN}_{\text{G,CWG}}(\cdot; \mathbf{W})$: A convolution implementation using the `WINOGRAD_NONFUSED` algorithm from the cuDNN library [7] on NVIDIA GPUs. It is based on the Winograd convolution algorithm [22], which runs faster but has higher numerical error compared to others.

For a given implementation $\text{NN}_{\text{impl}}(\cdot; \mathbf{W})$, our method finds pairs of $(\mathbf{x}_0, \mathbf{x}_{\text{adv}})$ represented as single precision floating point numbers such that

1. \mathbf{x}_0 and \mathbf{x}_{adv} are in the dynamic range of images:

$$\min \mathbf{x}_0 \geq 0, \min \mathbf{x}_{\text{adv}} \geq 0, \max \mathbf{x}_0 \leq 1, \text{ and } \max \mathbf{x}_{\text{adv}} \leq 1$$
2. \mathbf{x}_{adv} falls in the perturbation space of \mathbf{x}_0 : $\|\mathbf{x}_{\text{adv}} - \mathbf{x}_0\|_{\infty} \leq \epsilon$
3. The verifier claims that the robustness specification (2) holds for \mathbf{x}_0
4. The implementation falsifies the claim of the verifier:

$$L_{\text{CW}}(\text{NN}_{\text{impl}}(\mathbf{x}_{\text{adv}}; \mathbf{W}), t_0) < 0$$

Note that the first two conditions are accurately defined for any implementation compliant with the IEEE-754 standard [19], because the computation only involves element-wise subtraction and max-reduction that incur no accumulated error. The `Gurobi` solver used by `MIPVerify` operates with double precision internally. Therefore, to ensure that our adversarial examples satisfy the constraints considered by the solver, we also require that the first two conditions hold for $\mathbf{x}'_{\text{adv}} = \text{float64}(\mathbf{x}_{\text{adv}})$ and $\mathbf{x}'_0 = \text{float64}(\mathbf{x}_0)$ that are double precision representations of \mathbf{x}_{adv} and \mathbf{x}_0 .

4 Exploiting a complete verifier

We present two observations crucial to the exploitation to be described later.

Observation 1: Tiny perturbations on the network input result in random output perturbations. We select an image \mathbf{x} for which the verifier claims that the network makes robust predictions. We plot $\|\text{NN}(\mathbf{x} + \delta; \mathbf{W}) - \text{NN}(\mathbf{x}; \mathbf{W})\|_{\infty}$ against $-10^{-6} \leq \delta \leq 10^{-6}$, where the addition of $\mathbf{x} + \delta$ is only applied on the single input element that has the largest gradient magnitude. As shown in Figure 1, the change of the output is highly nonlinear with respect to the change of the input, and a small perturbation could result in a large fluctuation. Note that the output fluctuation is caused by accumulated floating point error instead of nonlinearities in the network because pre-activation values of all the ReLU units have the same signs for both \mathbf{x} and $\mathbf{x} + \delta$.

Observation 2: Different neural network implementations exhibit different floating point error characteristics. We evaluate the implementations on the whole MNIST test set and compare the outputs of the first layer (i.e., with only one linear transformation applied to the input) against that of $\text{NN}_{\text{C,M}}$. Figure 2 presents the histogram which shows that different implementations usually manifest different error behavior.

Method Overview: Given a network and weights $\text{NN}(\cdot; \mathbf{W})$, we search for image pairs $(\mathbf{x}_0, \mathbf{x}_1)$ such that the network is verifiably robust with respect to \mathbf{x}_0 , while $\mathbf{x}_1 \in \text{Adv}_{\epsilon}(\mathbf{x}_0)$ and $L_{\text{CW}}(\text{NN}(\mathbf{x}_1; \mathbf{W}), t_0)$ is less than the numerical fluctuation introduced by tiny input perturbations. We call \mathbf{x}_0 a *quasi-safe*

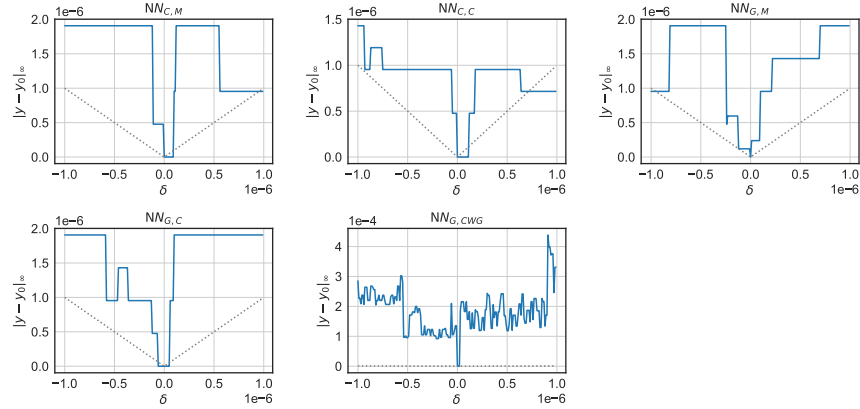


Fig. 1: Change of logits vector due to small single-element input perturbations for different implementations. The dashed lines are $y = |\delta|$.

image and \mathbf{x}_1 the corresponding *quasi-adversarial image*. Observation 1 suggests that an adversarial image might be obtained by randomly disturbing the quasi-adversarial image in the perturbation space. Observation 2 suggests that each implementation has its own adversarial images and needs to be handled separately. We search for the quasi-safe image by modifying the brightness of a natural image while querying a complete verifier whether it is near the boundary of robust predictions. Figure 3 illustrates this process.

Before explaining the details of our method, we first present the following proposition that formally establishes the existence of quasi-safe and quasi-adversarial images for continuous neural networks:

Proposition 1. *Let $E > 0$ be an arbitrarily small positive number. If a continuous neural network $\text{NN}(\cdot; \mathbf{W})$ can produce a robust classification for some input belonging to class t , and it does not constantly classify all inputs as class t , then there exists an input \mathbf{x}_0 such that*

$$0 < \min_{\mathbf{x} \in \text{Adv}_\epsilon(\mathbf{x}_0)} L_{\text{CW}}(\text{NN}(\mathbf{x}; \mathbf{W}), t) < E$$

Let $\mathbf{x}_1 = \text{argmin}_{\mathbf{x} \in \text{Adv}_\epsilon(\mathbf{x}_0)} L_{\text{CW}}(\text{NN}(\mathbf{x}; \mathbf{W}), t)$ be the minimizer of the above function. We call \mathbf{x}_0 a *quasi-safe image* and \mathbf{x}_1 a *quasi-adversarial image*.

Proof. Let $f(\mathbf{x}) := \min_{\mathbf{x}' \in \text{Adv}_\epsilon(\mathbf{x})} L_{\text{CW}}(\text{NN}(\mathbf{x}'; \mathbf{W}), t)$. Since $f(\cdot)$ is composed of continuous functions, $f(\cdot)$ is continuous. Suppose $\text{NN}(\cdot; \mathbf{W})$ is robust with respect to \mathbf{x}_+ that belongs to class t . Let \mathbf{x}_- be any input such that $L_{\text{CW}}(\text{NN}(\mathbf{x}_-; \mathbf{W}), t) < 0$, which exists because $\text{NN}(\cdot; \mathbf{W})$ does not constantly classify all inputs as class t . We have $f(\mathbf{x}_+) > 0$ and $f(\mathbf{x}_-) < 0$. The Poincaré-Miranda theorem asserts the existence of \mathbf{x}_0 such that $0 < f(\mathbf{x}_0) < E$.

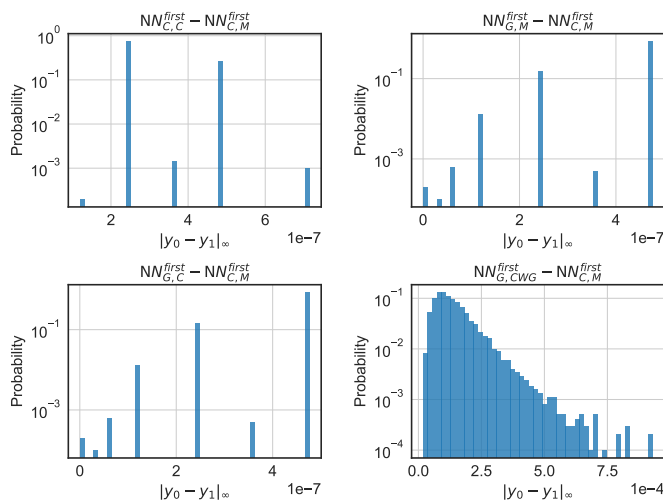


Fig. 2: Distribution of difference relative to $\text{NN}_{C,M}$ of first layer evaluated on MNIST test images.

Given a particular implementation $\text{NN}_{\text{impl}}(\cdot; \mathbf{W})$ and a natural image \mathbf{x}_{seed} which the network robustly classifies as class t_0 according to the verifier, we construct an adversarial input pair $(\mathbf{x}_0, \mathbf{x}_{\text{adv}})$ that meets the constraints described in Section 3 in three steps:

Step 1: We search for a coefficient $\alpha \in [0, 1]$ such that $\mathbf{x}_0 = \alpha \mathbf{x}_{\text{seed}}$ serves as the quasi-safe image. Specifically, we require the verifier to claim that the network is robust for $\alpha \mathbf{x}_{\text{seed}}$ but not so for $\alpha' \mathbf{x}_{\text{seed}}$ with $0 < (\delta = \alpha - \alpha') < \epsilon_r$, where ϵ_r should be small enough to allow quasi-adversarial images sufficiently close to the boundary. We set $\epsilon_r = 10^{-7}$. We use binary search to minimize δ starting from $\alpha' \leftarrow 0, \alpha \leftarrow 1$. We found that the MILP solver often becomes extremely slow when δ is small, so we start with binary search and switch to grid search by dividing the best known δ to 16 intervals if the solver exceeds a time limit.

Step 2: We search for the quasi-adversarial image \mathbf{x}_1 corresponding to \mathbf{x}_0 . We define a loss function with a tolerance of τ as $L(\mathbf{x}, \tau) := L_{\text{CW}}(\text{NN}(\mathbf{x}; \mathbf{W}), t_0) - \tau$, which can be incorporated in any verifier by modifying the bias of the Softmax layer. We aim to find τ_0 and τ_1 , where τ_0 is the minimal confidence of all images in the perturbation space of \mathbf{x}_0 , and τ_1 is slightly larger than τ_0 with \mathbf{x}_1 being the corresponding adversarial image. Formally:

$$\begin{cases} \forall \mathbf{x} \in \text{Adv}_\epsilon(\mathbf{x}_0) : L(\mathbf{x}, \tau_0) > 0 \\ \mathbf{x}_1 \in \text{Adv}_\epsilon(\mathbf{x}_0) \\ L(\mathbf{x}_1, \tau_1) < 0 \\ \tau_1 - \tau_0 < 10^{-7} \end{cases}$$

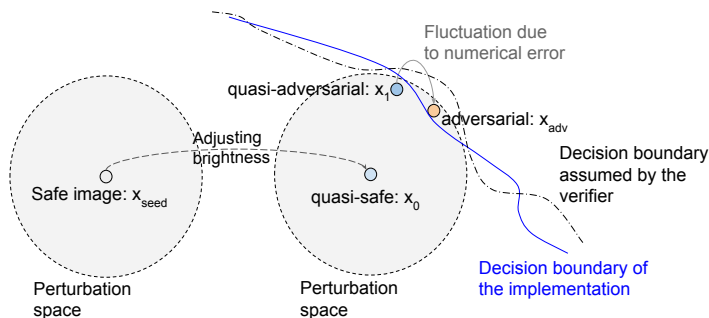


Fig. 3: Illustration of our method. Since the verifier does not model the floating point arithmetic details of the implementation, their decision boundaries for the classification problem diverge, which allows us to find adversarial inputs that cross the boundary via numerical error fluctuations. Note that the verifier usually does not comply with a well defined specification of $\text{NN}(\cdot; \mathbf{W})$, and therefore it does not define a decision boundary. The dashed boundary in the diagram is just for illustrative purposes.

Note that \mathbf{x}_1 is produced by the complete verifier as proof of nonrobustness given the tolerance τ_1 . The above values are found via binary search with initialization $\tau_0 \leftarrow 0$ and $\tau_1 \leftarrow L_{CW}(\text{NN}(\mathbf{x}_0; \mathbf{W}), t_0)$. In addition, we accelerate the binary search if the verifier can compute the *worst* objective defined as:

$$\tau_w = \min_{\mathbf{x} \in \text{Adv}_\epsilon(\mathbf{x}_0)} L_{CW}(\text{NN}(\mathbf{x}; \mathbf{W}), t_0) \quad (4)$$

In this case, we initialize $\tau_0 \leftarrow \tau_w - \delta_s$ and $\tau_1 \leftarrow \tau_w + \delta_s$. We empirically set $\delta_s = 3 \times 10^{-6}$ to incorporate the numerical error in the verifier so that $L(\mathbf{x}_0, \tau_w - \delta_s) > 0$ and $L(\mathbf{x}_0, \tau_w + \delta_s) < 0$. The binary search is aborted if the solver times out.

Step 3: We minimize $L_{CW}(\text{NN}(\mathbf{x}_1; \mathbf{W}), t_0)$ with hill climbing via applying small random perturbations on the quasi-adversarial image \mathbf{x}_1 while projecting back to $\text{Adv}_\epsilon(\mathbf{x}_0)$ to find an adversarial example. The perturbations are applied on patches of \mathbf{x}_1 as described in Algorithm 1.

Algorithm 1 Searching adversarial examples via hill climbing

- Input:** quasi-safe image \mathbf{x}_0
- Input:** target class number t
- Input:** quasi-adversarial image \mathbf{x}_1
- Input:** input perturbation bound ϵ
- Input:** a neural network inference implementation $\text{NN}_{\text{impl}}(\cdot; \mathbf{W})$
- Input:** number of iterations N (default value 1000)
- Input:** perturbation scale u (default value $2e-7$)
- Output:** an adversarial image \mathbf{x}_{adv} or FAILED


```

for Index  $i$  of  $\mathbf{x}_0$  do           ▷ Find the bounds  $\mathbf{x}_l$  and  $\mathbf{x}_u$  for allowed perturbations
   $x_l[i] \leftarrow \max(\text{nextafter}(x_0[i] - \epsilon, 0), 0)$ 
   $x_u[i] \leftarrow \min(\text{nextafter}(x_0[i] + \epsilon, 1), 1)$ 
  while  $x_0[i] - x_l[i] > \epsilon$  or  $\text{float64}(x_0[i]) - \text{float64}(x_l[i]) > \epsilon$  do
     $x_l[i] \leftarrow \text{nextafter}(x_l[i], 1)$ 
  end while
  while  $x_u[i] - x_0[i] > \epsilon$  or  $\text{float64}(x_u[i]) - \text{float64}(x_0[i]) > \epsilon$  do
     $x_u[i] \leftarrow \text{nextafter}(x_u[i], 0)$ 
  end while
end for

  ▷ We select the offset and stride based on the implementation to ensure that
  perturbed tiles contribute independently to the output. The Winograd algorithm in
  cuDNN produces  $9 \times 9$  output tiles for  $13 \times 13$  input tiles and  $5 \times 5$  kernels.
if  $\text{NN}_{\text{impl}}(\cdot; \mathbf{W})$  is  $\text{NN}_{\text{G,CWG}}(\cdot; \mathbf{W})$  then (offset, stride)  $\leftarrow (4, 9)$ 
else (offset, stride)  $\leftarrow (0, 4)$ 
end if

for  $i \leftarrow 1$  to  $N$  do
  for  $(h, w) \leftarrow (0, 0)$  to  $(\text{height}(\mathbf{x}_1), \text{width}(\mathbf{x}_1))$  step (stride, stride) do
     $\delta \leftarrow \text{uniform}(-u, u, (\text{stride} - \text{offset}, \text{stride} - \text{offset}))$ 
     $\mathbf{x}'_1 \leftarrow \mathbf{x}_1[:]$ 
     $\mathbf{x}'_1[h + \text{offset} : h + \text{stride}, w + \text{offset} : w + \text{stride}] += \delta$ 
     $\mathbf{x}'_1 \leftarrow \max(\min(\mathbf{x}'_1, \mathbf{x}_u), \mathbf{x}_l)$ 
    if  $L_{\text{CW}}(\text{NN}_{\text{impl}}(\mathbf{x}'_1; \mathbf{W}), t) < L_{\text{CW}}(\text{NN}_{\text{impl}}(\mathbf{x}_1; \mathbf{W}), t)$  then
       $\mathbf{x}_1 \leftarrow \mathbf{x}'_1$ 
    end if
  end for
end for
if  $L_{\text{CW}}(\text{NN}_{\text{impl}}(\mathbf{x}_1; \mathbf{W}), t) < 0$  then return  $\mathbf{x}_{\text{adv}} \leftarrow \mathbf{x}_1$ 
else return FAILED
end if

```

Experiments: We conduct our experiments on a workstation with an NVIDIA Titan RTX GPU and an AMD Ryzen Threadripper 2970WX CPU. We train the small architecture from Xiao et al. [48] with the PGD adversary and the RS Loss on MNIST and CIFAR10 datasets. The network has two convolutional layers with 4×4 filters, 2×2 stride, and 16 and 32 output channels, respectively, and two fully connected layers with 100 and 10 output neurons. The trained networks achieve 94.63% and 44.73% provable robustness with perturbations of ℓ_∞ bounded by 0.1 and $2/255$ on the two datasets, respectively, similar to the results reported in Xiao et al. [48]. Our code is available at <https://github.com/jia-kai/realadv>.

Although our method only needs $O(-\log \epsilon)$ invocations of the verifier where ϵ is the threshold in the binary search, the verifier still takes most of the time and is too slow for a large benchmark. Therefore, for each dataset, we test our method on 32 images randomly sampled from the verifiably robustly classified test images. All the implementations that we have considered are successfully exploited. Specifically, our benchmark contains $32 \times 2 \times 5 = 320$ cases,

Table 1: Number of adversarial examples successfully found for different neural network inference implementations

	$NN_{C,M}$	$NN_{C,C}$	$NN_{G,M}$	$NN_{G,C}$	$NN_{G,CWG}$
MNIST	2	3	1	3	7
CIFAR10	16	12	7	6	25

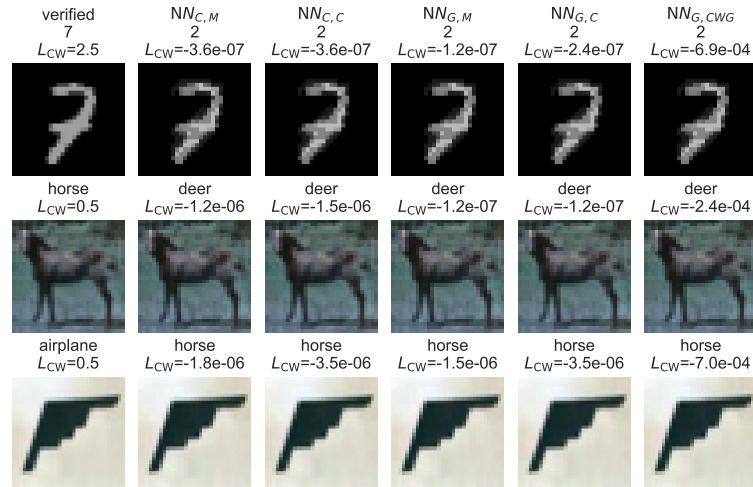


Fig. 4: The quasi-safe images with respect to which all implementations are successfully exploited, and the corresponding adversarial images.

while adversarial examples are found for 82 of them. The failed cases correspond to large τ_1 values in Step 2 due to verifier timeouts or the discrepancy of floating point arithmetic between the verifier and the implementations. Let $\tau_{\text{impl}} := L_{CW}(\text{NN}_{\text{impl}}(\mathbf{x}_1; \mathbf{W}), t)$ denote the loss of a quasi-adversarial input on a particular implementation. Algorithm 1 succeeds on all cases with $\tau_{\text{impl}} < 8.3 \times 10^{-7}$ (35 such cases in total), while 18 among them have $\tau_{\text{impl}} < 0$ due to floating point discrepancy (i.e., the quasi-adversarial input is already an adversarial input for this implementation). The most challenging case (i.e., with largest τ_{impl}) on which Algorithm 1 succeeds has $\tau_{\text{impl}} = 3.2 \times 10^{-4}$. The largest value of τ_{impl} is 3.7. Table 1 presents the detailed numbers for each implementation. Figure 4 shows the quasi-safe images on which our exploitation method succeeds for all implementations and the corresponding adversarial images.

5 Exploiting an incomplete verifier

The relaxation adopted in certification methods renders them incomplete but also makes their verification claims more robust to floating point error compared

to complete verifiers. In particular, we evaluate the CROWN framework [49] on our randomly selected MNIST test images and the corresponding quasi-safe images from Section 4. CROWN is able to verify the robustness of the network on 29 out of the 32 original test images, but it is unable to prove the robustness for any of the quasi-safe images. Note that MIPVerify claims that the network is robust with respect to all the original test images and the corresponding quasi-safe images.

Incomplete verifiers are still vulnerable if we allow arbitrary network architectures and weights. Our exploitation builds on the observation that verifiers typically need to merge always-active ReLU units with their subsequent layers to reduce the number of nonlinearities and achieve a reasonable speed. The merge of layers involves computing merged “equivalent” weights, which is different from the floating point computation adopted by an inference implementation.

We build a neural network that takes a 13×13 single-channel input image, followed by a 5×5 convolutional layer with a single output channel, two fully connected layers with 16 output neurons each, a fully connected layer with one output neuron denoted as $u = \max(\mathbf{W}_u \mathbf{h}_u + b_u, 0)$, and a final linear layer that computes $\mathbf{y} = [u; 10^{-7}]$ as the logits vector. All the hidden layers have ReLU activation. The input \mathbf{x}_0 is taken from a Gaussian distribution. The hidden layers have random Gaussian coefficients, and the biases are chosen so that (i) the ReLU neurons before u are always activated for inputs in the perturbation space of \mathbf{x}_0 , and (ii) the neuron u is never activated while b_u is the maximum possible value (i.e., $b_u = -\max_{\mathbf{x} \in \text{Adv}_\epsilon(\mathbf{x}_0)} \mathbf{W}_u \mathbf{h}_u(\mathbf{x})$). CROWN is able to prove that all ReLU neurons before u are always activated but u is never activated, and therefore it claims that the network is robust with respect to perturbations around \mathbf{x}_0 . However, by initializing the quasi-adversarial input $\mathbf{x}_1 \leftarrow \mathbf{x}_0 + \epsilon \text{sign}(\mathbf{W}_{equiv})$ where \mathbf{W}_{equiv} is the product of all the coefficient matrices of the layers up to u , we successfully find adversarial inputs for all the five implementations considered in this work by randomly perturbing \mathbf{x}_1 using Algorithm 1 with a larger number of iterations ($N = 10000$) due to the smaller input size.

Note that the output scores can be manipulated to appear less suspicious. For instance, we can set $\mathbf{z} = \text{clip}(10^7 \cdot \mathbf{y}, -2, 2)$ as the final output in the above example so that \mathbf{z} becomes a more “naturally looking” classification score in the range $[-2, 2]$ and its perturbation due to floating point error is also enlarged to the unit scale. The extreme constants 10^{-7} and 10^7 can also be obfuscated by using multiple consecutive scaling layers with each one having small scaling factors such as 0.1 and 10.

6 Discussion

We have shown that some neural network verifiers are systematically exploitable. One appealing remedy is to introduce relaxations into complete verifiers, such as by verifying for a larger ϵ or setting a threshold for accepted confidence score. For example, it might be tempting to claim the robustness of a network for $\epsilon = 0.09999$ if it is verified for $\epsilon = 0.1$. We emphasize that there are no guarantees provided by any floating point complete verifier currently extant. Moreover, the

difference between the true robust perturbation bound and the bound claimed by an unsound verifier might be much larger if the network has certain properties. For example, `MIPVerify` has been observed to give NaN results when verifying pruned neural networks [15]. The adversary might also be able to manipulate the network to scale the scores arbitrarily, as discussed in Section 5. The correct solution requires obtaining a tight relaxation bound that is sound for both the verifier and the inference implementation, which is extremely challenging.

A possible fix for complete verification is to adopt exact MILP solvers with rational inputs [39]. There are three challenges: (i) The efficiency of exactly solving the large amounts of computation in neural network inference has not been studied and is unlikely to be satisfactory, (ii) The computation that derives the MILP formulation from a verification specification, such as the neuron bound analysis in Tjeng et al. [42], must also be exact, but existing neural network verifiers have not attempted to define and implement exact arithmetic with the floating point weights, and (iii) The results of exact MILP solvers are only valid for an exact neural network inference implementation, but such exact implementations are not widely available (not provided by any deep learning libraries that we are aware of), and their efficiency remains to be studied.

Alternatively, one may obtain sound and nearly complete verification by adopting a conservative MILP solver based on techniques such as directed rounding [28]. We also need to ensure all arithmetic in the verifier to derive the MILP formulation soundly over-approximates floating point error. This is more computationally feasible than exact verification discussed above. It is similar to the approach used in some sound incomplete verifiers that incorporate floating point error by maintaining upper and lower rounding bounds of internal computations [36, 37]. However, this approach relies on the specific implementation details of the inference algorithm — optimizations such as Winograd [22] or FFT [1], or deployment in hardware accelerators with lower floating point precision such as Bfloat16 [4], would either invalidate the robustness guarantees or require changes to the analysis algorithm. Therefore, we suggest that these sound verifiers explicitly state the requirements on the inference implementations for which their results are sound. A possible future research direction is to devise a universal sound verification framework that can incorporate different inference implementations.

Another approach for sound and complete neural network verification is to quantize the computation to align the inference implementation with the verifier. For example, if we require all activations to be multiples of s_0 and all weights to be multiples of s_1 , where $s_0s_1 > 2E$ and E is a very loose bound of possible implementation error, then the output can be rounded to multiples of s_0s_1 to completely eliminate numerical error. Binarized neural networks [18] are a family of extremely quantized networks, and their verification [20, 27, 35] is sound and complete. However, the problem of robust training and verification of quantized neural networks [20] is relatively under-examined compared to that of real-valued neural networks [24, 26, 42, 48].

7 Conclusion

Floating point error should not be overlooked in the verification of real-valued neural networks, as we have presented techniques that efficiently find witnesses for the unsoundness of two verifiers. Unfortunately, floating point soundness issues have not received sufficient attention in neural network verification research. A user has few choices if they want to obtain sound verification results for a neural network, especially if they deploy accelerated neural network inference implementations. We hope our results will help to guide future neural network verification research by providing another perspective on the tradeoff between soundness, completeness, and scalability.

Acknowledgments We would like to thank Gagandeep Singh and Kai Xiao for providing invaluable suggestions on an early manuscript.

References

1. Abtahi, T., Shea, C., Kulkarni, A., Mohsenin, T.: Accelerating convolutional neural network with FFT on embedded hardware. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* **26**(9), 1737–1749 (2018)
2. Boldo, S., Melquiond, G.: *Computer Arithmetic and Formal Proofs: Verifying Floating-point Algorithms with the Coq System*. Elsevier (2017)
3. Bunel, R., Lu, J., Turkaslan, I., Kohli, P., Torr, P., Mudigonda, P.: Branch and bound for piecewise linear neural network verification. *Journal of Machine Learning Research* **21**(2020) (2020)
4. Burgess, N., Milanovic, J., Stephens, N., Monachopoulos, K., Mansell, D.: Bfloat16 processing for neural networks. In: 2019 IEEE 26th Symposium on Computer Arithmetic (ARITH), pp. 88–91, IEEE (2019)
5. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: 2017 IEEE Symposium on Security and Privacy (SP), pp. 39–57, IEEE (2017)
6. Cheng, C.H., Nührenberg, G., Ruess, H.: Maximum resilience of artificial neural networks. In: *International Symposium on Automated Technology for Verification and Analysis*, pp. 251–268, Springer (2017)
7. Chetlur, S., Woolley, C., Vandermersch, P., Cohen, J., Tran, J., Catanzaro, B., Shelhamer, E.: cudnn: Efficient primitives for deep learning. *arXiv preprint arXiv:1410.0759* (2014)
8. Corzilius, F., Loup, U., Junges, S., Ábrahám, E.: Smt-rat: an smt-compliant non-linear real arithmetic toolbox. In: *International Conference on Theory and Applications of Satisfiability Testing*, pp. 442–448, Springer (2012)
9. Das, A., Briggs, I., Gopalakrishnan, G., Krishnamoorthy, S., Panckekha, P.: Scalable yet rigorous floating-point error analysis. In: *SC20*, pp. 1–14, IEEE (2020)
10. Dutta, S., Jha, S., Sankaranarayanan, S., Tiwari, A.: Output range analysis for deep feedforward neural networks. In: *NASA Formal Methods Symposium*, pp. 121–138, Springer (2018)
11. Dvijotham, K., Gowal, S., Stanforth, R., Arandjelovic, R., O’Donoghue, B., Uesato, J., Kohli, P.: Training verified learners with learned verifiers. *arXiv preprint arXiv:1805.10265* (2018)

12. Ehlers, R.: Formal verification of piece-wise linear feed-forward neural networks. In: International Symposium on Automated Technology for Verification and Analysis, pp. 269–286, Springer (2017)
13. Fischetti, M., Jo, J.: Deep neural networks and mixed integer linear optimization. *Constraints* **23**(3), 296–309 (2018)
14. 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: 2018 IEEE Symposium on Security and Privacy (SP), pp. 3–18, IEEE (2018)
15. Guidotti, D., Leofante, F., Pulina, L., Tacchella, A.: Verification of neural networks: Enhancing scalability through pruning. arXiv preprint arXiv:2003.07636 (2020)
16. Gurobi Optimization, L.: Gurobi optimizer reference manual (2020), URL <http://www.gurobi.com>
17. Huang, X., Kwiatkowska, M., Wang, S., Wu, M.: Safety verification of deep neural networks. In: International Conference on Computer Aided Verification, pp. 3–29, Springer (2017)
18. Hubara, I., Courbariaux, M., Soudry, D., El-Yaniv, R., Bengio, Y.: Binarized neural networks. In: NeurIPS, pp. 4107–4115, Curran Associates, Inc. (2016)
19. IEEE: IEEE standard for floating-point arithmetic. IEEE Std 754-2008 pp. 1–70 (2008)
20. Jia, K., Rinard, M.: Efficient exact verification of binarized neural networks. In: NeurIPS, vol. 33, pp. 1782–1795, Curran Associates, Inc. (2020)
21. Katz, G., Barrett, C., Dill, D.L., Julian, K., Kochenderfer, M.J.: Reluplex: An efficient smt solver for verifying deep neural networks. In: International Conference on Computer Aided Verification, pp. 97–117, Springer (2017)
22. Lavin, A., Gray, S.: Fast algorithms for convolutional neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4013–4021 (2016)
23. Lomuscio, A., Maganti, L.: An approach to reachability analysis for feed-forward relu neural networks. arXiv preprint arXiv:1706.07351 (2017)
24. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. In: ICLR (2018)
25. Miné, A.: Relational abstract domains for the detection of floating-point run-time errors. In: European Symposium on Programming, pp. 3–17, Springer (2004)
26. Mirman, M., Gehr, T., Vechev, M.: Differentiable abstract interpretation for provably robust neural networks. In: Dy, J., Krause, A. (eds.) Proceedings of the 35th International Conference on Machine Learning, Proceedings of Machine Learning Research, vol. 80, pp. 3578–3586, PMLR, Stockholmsmässan, Stockholm Sweden (10–15 Jul 2018)
27. Narodytska, N., Kasiviswanathan, S., Ryzhyk, L., Sagiv, M., Walsh, T.: Verifying properties of binarized deep neural networks. In: Thirty-Second AAAI Conference on Artificial Intelligence (2018)
28. Neumaier, A., Shcherbina, O.: Safe bounds in linear and mixed-integer linear programming. *Mathematical Programming* **99**(2), 283–296 (2004)
29. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: PyTorch: An imperative style, high-performance deep learning library. In: NeurIPS, pp. 8024–8035, Curran Associates, Inc. (2019)
30. Raghu, M., Schmidt, E.: A survey of deep learning for scientific discovery. *ArXiv abs/2003.11755* (2020)

31. Raghunathan, A., Steinhardt, J., Liang, P.S.: Semidefinite relaxations for certifying robustness to adversarial examples. In: NeurIPS, pp. 10877–10887, Curran Associates, Inc. (2018)
32. Rümmer, P., Wahl, T.: An smt-lib theory of binary floating-point arithmetic. In: International Workshop on Satisfiability Modulo Theories (SMT), p. 151 (2010)
33. Salman, H., Yang, G., Zhang, H., Hsieh, C.J., Zhang, P.: A convex relaxation barrier to tight robustness verification of neural networks. In: NeurIPS, pp. 9832–9842 (2019)
34. Scheibler, K., Winterer, L., Wimmer, R., Becker, B.: Towards verification of artificial neural networks. In: MBMV, pp. 30–40 (2015)
35. Shih, A., Darwiche, A., Choi, A.: Verifying binarized neural networks by anguine-style learning. In: International Conference on Theory and Applications of Satisfiability Testing, pp. 354–370, Springer (2019)
36. Singh, G., Gehr, T., Mirman, M., Püschel, M., Vechev, M.: Fast and effective robustness certification. In: NeurIPS, pp. 10802–10813, Curran Associates, Inc. (2018)
37. Singh, G., Gehr, T., Püschel, M., Vechev, M.T.: An abstract domain for certifying neural networks. *Proceedings of the ACM on Programming Languages* **3**, 1 – 30 (2019)
38. Solovyev, A., Baranowski, M.S., Briggs, I., Jacobsen, C., Rakamarić, Z., Gopalakrishnan, G.: Rigorous estimation of floating-point round-off errors with symbolic taylor expansions. *TOPLAS* **41**(1), 1–39 (2018)
39. Steffy, D.E., Wolter, K.: Valid linear programming bounds for exact mixed-integer programming. *INFORMS Journal on Computing* **25**(2), 271–284 (2013)
40. Szegedy, C., Zaremba, W., Sutskever, I., Estrach, J.B., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. In: ICLR (2014)
41. Titolo, L., Feliú, M.A., Moscato, M.M., Muñoz, C.A.: An abstract interpretation framework for the round-off error analysis of floating-point programs. In: VMCAI, pp. 516–537 (2018)
42. Tjeng, V., Xiao, K.Y., Tedrake, R.: Evaluating robustness of neural networks with mixed integer programming. In: ICLR (2019)
43. Tramer, F., Boneh, D.: Adversarial training and robustness for multiple perturbations. In: NeurIPS, pp. 5866–5876, Curran Associates, Inc. (2019)
44. Wang, S., Pei, K., Whitehouse, J., Yang, J., Jana, S.: Formal security analysis of neural networks using symbolic intervals. In: 27th USENIX Security Symposium (USENIX Security 18), pp. 1599–1614 (2018)
45. Weng, L., Zhang, H., Chen, H., Song, Z., Hsieh, C.J., Daniel, L., Boning, D., Dhillon, I.: Towards fast computation of certified robustness for ReLU networks. In: International Conference on Machine Learning, pp. 5276–5285 (2018)
46. Wong, E., Kolter, J.Z.: Provable defenses against adversarial examples via the convex outer adversarial polytope. arXiv preprint arXiv:1711.00851 (2017)
47. Wong, E., Rice, L., Kolter, J.Z.: Fast is better than free: Revisiting adversarial training. In: ICLR (2020)
48. Xiao, K.Y., Tjeng, V., Shafiqullah, N.M.M., Madry, A.: Training for faster adversarial robustness verification via inducing reLU stability. In: ICLR (2019)
49. Zhang, H., Weng, T.W., Chen, P.Y., Hsieh, C.J., Daniel, L.: Efficient neural network robustness certification with general activation functions. In: NeurIPS, pp. 4939–4948, Curran Associates, Inc. (2018)