

# Minimum Uncertainty Based Detection of Adversaries in Deep Neural Networks

Fatemeh Sheikholeslami, *Student Member, IEEE*, Swayambhoo Jain, *Member, IEEE*,  
and Georgios B. Giannakis, *Fellow, IEEE*

**Abstract**—Despite their unprecedented performance in various domains, utilization of Deep Neural Networks (DNNs) in safety-critical environments is severely limited in the presence of even small adversarial perturbations. The present work develops a randomized approach to detecting such perturbations based on minimum uncertainty metrics that rely on sampling at the hidden layers during the DNN inference stage. Inspired by Bayesian approaches to uncertainty estimation, the sampling probabilities are designed for effective detection of the adversarially corrupted inputs. Being modular, the novel detector of adversaries can be conveniently employed by any pre-trained DNN at no extra training overhead. Selecting which units to sample per hidden layer entails quantifying the amount of DNN output uncertainty, where the overall uncertainty is expressed in terms of its layer-wise components - what also promotes scalability. Sampling probabilities are then sought by minimizing uncertainty measures layer-by-layer, leading to a novel convex optimization problem that admits an exact solver with superlinear convergence rate. By simplifying the objective function, low-complexity approximate solvers are also developed. In addition to valuable insights, these approximations link the novel approach with state-of-the-art randomized adversarial detectors. The effectiveness of the novel detectors in the context of competing alternatives is highlighted through extensive tests for various types of adversarial attacks with variable levels of strength.

**Index Terms**—Adversarial input, Bayesian neural networks, attack detection, uncertainty estimation.

## 1 INTRODUCTION

Unprecedented learning capability offered by Deep Neural Networks (DNNs) has enabled state-of-the-art performance in diverse tasks such as object recognition and detection [1], [2], [3], speech recognition and language translation [4], voice synthesis [5], and many more, to reach or even surpass human-level accuracy. Despite their performance however, recent studies have cast doubt on the reliability of DNNs as highly-accurate networks are shown to be extremely sensitive to carefully crafted inputs designed to fool them [6], [7], [8]. Such fragility can easily lead to sabotage once adversarial entities target critical environments such as autonomous cars [9], automatic speech recognition [10], and face detection [3], [11], [12]. The extreme brittleness of convolutional neural networks (CNNs) for image classification is highlighted since small adversarial perturbations on the clean image, although often imperceptible to the human eye, can lead the trained CNNs to classify the *adversarial examples* incorrectly with high confidence. In particular, design of powerful adversarial perturbations in environments with different levels of complexity and knowledge about the target CNN, known as white, grey, and black-box attacks, have been investigated in several works [7], [13], [14], [15], [16], [17]. These considerations motivate well the need for designing robust and powerful attack detection mechanisms for reliable and safe utilization of DNNs [18].

Defense against adversarial perturbations has been mainly pursued in two broad directions: (i) attack *detection*,

and (ii) attack *recovery*. Methods in the first category aim at detecting adversarial corruption in the input by classifying the input as clean or adversarial, based on tools as diverse as auto-encoders [19], detection sub-networks [20], [21], and dropout units [22]. On the other hand, methods in the second category are based on recovery schemes that robustify the classification by data pre-processing and randomization [23], [24], [25], adversarial training [26], [27], [28], [29], [30], sparsification of the network [31], [32] and Lipschitz regularization [33], [34], to name just a few.

Furthermore, the so-termed *over-confidence* of DNNs in classifying “out-of-distribution,” meaning samples which lie in unexplored regions of the input domain, or even “misclassified” samples, has been unraveled in [35], [36]. This has motivated the need for uncertainty estimation as well as calibration of the networks for robust classification. Modern Bayesian neural networks target this issue by modeling the distribution of DNN weights as random [37], and estimating the DNN output uncertainty through predictive entropy, variance, or mutual information [22], [38], [39], [40]. The well-known dropout regularization technique is one such approximate Bayesian neural network, now widely used in training and testing of DNNs [41], [42], [43].

Moreover, approaches relying on dropout units have shown promising performance in successfully detecting adversarial attacks, where other defense mechanisms fail [13]. In particular, [22] utilizes randomness of dropout units during the test phase as a defense mechanism, and approximates the classification uncertainty by Monte Carlo (MC) estimation of the output variance. Based on the latter, images with high classification uncertainty are declared as adversarial. Recently, dropout defense has been generalized to non-uniform sampling [44], where entries of the hidden-layers

• Part of this work was done during a summer research internship at Technicolor AI Lab in Palo Alto, CA - USA. This research was supported in part by NSF grant 151405\6, 1505970, 1901134, and 1711471. Author emails: sheik081@umn.edu, swayambhoo.jain@gmail.com, georgios@umn.edu

are randomly sampled, with probabilities proportional to the entry values. This heuristic sampling of units per layer is inspired by intuitive reasoning: activation units with large entries have more information and should be sampled more often [44]. However, analytical understanding has not been investigated.

The goal here is to further expand the understanding of uncertainty estimation in DNNs, and thereby improve the detection of adversarial inputs. The premise is that inherent distance of the adversarial perturbation from the natural-image manifold will cause the overall network uncertainty to exceed that of the clean image, and thus successful detection can be obtained.

To this end, and inspired by [44], we rely on random sampling of units per hidden layer of a pre-trained network to introduce randomness. Moreover, inspired by the Bayesian approaches to uncertainty estimation, the overall uncertainty of a given image is then quantified in terms of its hidden-layer components. We then formulate the task of adversary detection as uncertainty minimization by optimizing over the sampling probabilities to provide effective detection. Subsequently, we develop an exact solver with super-linear convergence rate as well as approximate low-complexity solvers for an efficient layer-by-layer uncertainty minimization scheme. Furthermore, we draw connections with uniform dropout [22] as well as stochastic approximate pruning (SAP) [44], and provide an efficient implementation of the novel approach by interpreting it as a non-uniform dropout scheme. Extensive numerical tests on CIFAR10 and high-quality cats-and-dogs images in the presence of various attack schemes corroborate the importance of our designs of sampling probabilities, as well as the placement of sampling units per hidden layer for improved detection of adversarial inputs.

The rest of the paper is organized as follows. An overview on Bayesian inference and uncertainty-based detection in neural networks is provided in Section 2. Inspired by this, the proposed class of detectors is introduced in Section 3, and exact as well as low-complexity approximate solvers for the layer-by-layer uncertainty minimization are the subjects of Section 4. Implementation issues are dealt with in Section 5, numerical tests are provided in Section 6, and concluding remarks are discussed in Section 7.

## 2 BAYESIAN NEURAL NETWORK PRELIMINARIES

Bayesian inference is among the powerful tools utilized for analytically understanding and quantifying uncertainty in DNNs [42], [45]. In this section, we provide a short review on the basics of Bayesian neural networks, and move on to the inference phase for adversary detection in Section 2.2, which is of primary interest in this work.

Consider an  $L$ -layer deep neural network, which maps the input  $\mathbf{x} \in \mathcal{X}$  to output  $\mathbf{y} \in \mathcal{Y}$ . The weights are denoted by  $\omega := \{\mathbf{W}_l\}_{l=1}^L$ , and are modeled as random variables with prior probability density function (pdf)  $p(\omega)$ .

Given training input  $\mathbf{X} := [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$  and output data  $\mathbf{Y} := [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n]$ , it is assumed that the parameters  $\omega$  only depend on these  $(\mathbf{X}, \mathbf{Y})$  data. As a result, the

predictive pdf for a new input  $\mathbf{x}_\nu$  can be obtained via marginalization as [41]

$$p(\mathbf{y}_\nu | \mathbf{x}_\nu, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}_\nu | \mathbf{x}_\nu, \omega) p(\omega | \mathbf{X}, \mathbf{Y}) d\omega \quad (1)$$

which requires knowing the conditional  $p(\omega | \mathbf{X}, \mathbf{Y})$ . The complexity of estimating  $p(\omega | \mathbf{X}, \mathbf{Y})$  motivates well the variational inference (VI) approach, where  $p(\omega | \mathbf{X}, \mathbf{Y})$  is replaced by a surrogate pdf  $q_\theta(\omega)$  that is parameterized by  $\theta$ . For  $q_\theta(\omega)$ , it is desired to: (D1) approximate closely  $p(\omega | \mathbf{X}, \mathbf{Y})$ ; and, (D2) provide easy marginalization in (1) either in closed form or empirically. To meet (D1), the surrogate is chosen by minimizing the Kullback-Leibler (KL) divergence  $KL(p(\omega | \mathbf{X}, \mathbf{Y}), q_\theta(\omega))$ , which is subsequently approximated by the log evidence lower bound [46, p. 462]

$$\mathcal{L}_{VI}(\theta) := \int q_\theta(\omega) \log p(\mathbf{Y} | \mathbf{X}, \omega) d\omega - KL(q_\theta(\omega) || p(\omega)) . \quad (2)$$

Finding  $q_\theta$  boils down to maximizing the log evidence lower bound, that is,  $\theta_{VI} = \arg \max_\theta \mathcal{L}_{VI}(\theta)$ . A common choice for  $q_\theta(\omega)$  to also satisfy (D2) is described next.

### 2.1 Variational inference

A simple yet effective choice for  $q_\theta(\omega)$  is a factored form modeling the weights as independent across layers, that is

$$q_\theta(\omega) = \prod_{l=1}^L q(\mathbf{W}_l; \mathbf{M}_l, \boldsymbol{\theta}_{z_l}) \quad (3)$$

where the  $l$ -th layer with  $h_l$  hidden units is modeled as

$$\mathbf{W}_l = \mathbf{M}_l \text{diag}([z_{l,1}, z_{l,2}, \dots, z_{l,h_l}]) , \quad l = 1, \dots, L \quad (4)$$

where  $\mathbf{M}_l$  is an  $h_{l+1} \times h_l$  deterministic weight matrix multiplied by a diagonal matrix formed by the binary random vector  $\mathbf{z}_l := [z_{l,1}, z_{l,2}, \dots, z_{l,h_l}] \in \{0, 1\}^{h_l}$  with entries drawn from a pmf  $q_z(\mathbf{z}_l; \boldsymbol{\theta}_{z_l})$  parameterized by  $\boldsymbol{\theta}_{z_l}$ .

If the entries  $\{z_{l,i}\}$  are i.i.d. Bernoulli with (identical) probability (w.p.)  $\pi$ , they effect what is referred to as *uniform* (across layers and nodes) *dropout*, which is known to prevent overfitting [45]. Clearly, the parameter set  $\theta := \{\mathbf{M}_l, \boldsymbol{\theta}_{z_l}\}_{l=1}^L = \{\mathbf{M}_l\}_{l=1}^L \cup \{\pi\}$  fully characterizes  $q_\theta(\omega)$ . The dropout probability  $1 - \pi$  is preselected in practice, while  $\{\mathbf{M}_l\}_{l=1}^L$  can be obtained using the training data by maximizing the log evidence lower bound in (2). Nonetheless, integration in (2) over all the Bernoulli variables is analytically challenging, while sampling from the Bernoulli pmf is relatively cheap. This prompts approximate yet efficient integration using Monte Carlo estimation. A more detailed account of training Bayesian neural networks can be found in [41], [46], [47]. Moving on, the ensuing subsection deals with detection of adversarial inputs inspired by Bayesian neural networks.

### 2.2 Detection of DNN adversaries

In addition to facilitating ELBO maximization during the training phase, probabilistic view on the network parameters during the test phase can also be utilized towards output uncertainty estimation for the detection of adversarial inputs [41]. To do so, detection during the testing phase

proceeds by approximating the predictive pdf in (1) using the variational surrogate  $q_\theta(\omega)$ , as

$$p(\mathbf{y}_\nu | \mathbf{x}_\nu, \mathbf{X}, \mathbf{Y}) \approx \int p(\mathbf{y}_\nu | \mathbf{x}_\nu, \omega) q_\theta(\omega) d\omega. \quad (5)$$

Deciphering whether a given input  $\mathbf{x}_\nu$  is adversarial entails three steps: (S1) parametric modeling of  $q_\theta(\omega)$ ; (S2) estimating the DNN output *uncertainty* captured by  $p(\mathbf{y}_\nu | \mathbf{x}_\nu, \mathbf{X}, \mathbf{Y})$ ; and (S3) declaring  $\mathbf{x}_\nu$  as adversarial if the output uncertainty exceeds a certain threshold, and clean otherwise. These steps are elaborated next.

*Step 1: Parametric modeling of  $q_\theta(\omega)$ .* Recall that uniform dropout offers a popular special class of  $q_\theta(\omega)$  pdfs, and has been employed in adversary detection [22]. Here, we specify the richer model of  $q_\theta(\omega)$  in (3) and (4) that will turn out to markedly improve detection performance. Different from uniform dropout, we will allow for (possibly correlated) Bernoulli variables with carefully selected (possibly non-identical) parameters. If such general  $\{\boldsymbol{\theta}_{z_l}\}_{l=1}^L$  can be obtained, matrices  $\{\mathbf{M}_l\}_{l=1}^L$  are then found as follows.

Let  $\{\mathbf{W}_l^{(TR)}\}_{l=1}^L$  be deterministic weight matrices obtained via non-Bayesian training that we denote as (TR)<sup>1</sup>. We will use  $\mathbf{W}_l^{(TR)}$  to specify the mean of the random weight matrix  $\mathbf{W}_l$  in our approach, meaning we choose  $\mathbb{E}_{q_z(\mathbf{z}_l; \boldsymbol{\theta}_l)}[\mathbf{W}_l] \mathbf{x}_{(l-1)} = \mathbf{W}_l^{(TR)} \mathbf{x}_{(l-1)} \forall l$ , where  $\mathbf{x}_{(l-1)}$  is the output of the  $(l-1)$ st layer for a given input  $\mathbf{x}_\nu$  passing through the DNN with deterministic weights  $\{\mathbf{W}_l^{(TR)}\}_{l=1}^L$ . With  $\mathbf{W}_l^{(TR)}$  available, we first design  $q_z(\mathbf{z}_l; \boldsymbol{\theta}_{z_l})$ ; next, we find  $\boldsymbol{\theta}_{z_l}$ ; and then  $\mathbf{M}_l$ , as

$$\mathbf{M}_l = \mathbf{W}_l^{(TR)} \text{diag}^\dagger \left( \mathbb{E}_{q_z(\mathbf{z}_l; \boldsymbol{\theta}_{z_l})}[\mathbf{z}_l] \right), \quad l = 1, \dots, L \quad (6)$$

where the pseudo-inverse  $\dagger$  means that inverse entries are replaced with zeros if  $\mathbb{E}_{q_z(\mathbf{z}_l; \boldsymbol{\theta}_{z_l})}[z_{l,i}] = 0$ .

*Step 2: Quantifying the DNN output uncertainty.* Since evaluation of  $p(\mathbf{y}_\nu | \mathbf{x}_\nu, \mathbf{X}, \mathbf{Y})$  in (5) is prohibitive, one can estimate it using MC sampling. In particular, one can readily obtain MC estimates of (conditional) moments of  $\mathbf{y}_\nu$ . For instance, its mean and variance can be estimated as

$$\mathbb{E}_{q_\theta(\omega)}[\mathbf{y}_\nu | \mathbf{x}_\nu; \{\boldsymbol{\theta}_{z_l}\}_{l=1}^L] \simeq \bar{\mathbf{y}}_\nu = \frac{1}{R} \sum_{r=1}^R \mathbf{y}_\nu^{(r)}$$

and

$$\text{Cov}_{q_\theta(\omega)}[\mathbf{y}_\nu | \mathbf{x}_\nu; \{\boldsymbol{\theta}_{z_l}\}_{l=1}^L] \simeq \frac{1}{R} \sum_{r=1}^R \mathbf{y}_\nu^{(r)} \mathbf{y}_\nu^{(r)\top} - \bar{\mathbf{y}}_\nu \bar{\mathbf{y}}_\nu^\top \quad (7)$$

where  $\mathbf{y}_\nu^{(r)}$  is the output of the  $r$ -th DNN realized through weights  $\{\mathbf{W}_l^{(r)}\}_{l=1}^L$  with input  $\mathbf{x}_\nu$ . The predictive variance is the trace of  $\text{Cov}_{q_\theta(\omega)}[\mathbf{y}_\nu | \mathbf{x}_\nu; \{\boldsymbol{\theta}_{z_l}\}_{l=1}^L]$  that we henceforth abbreviate as  $\text{Tr}(\text{Cov}_{q_\theta(\omega)}[\mathbf{y}_\nu | \mathbf{x}_\nu])$ . Given  $\mathbf{x}_\nu$ , the latter has been used to quantify output uncertainty as  $U(\mathbf{x}_\nu) = \text{Tr}(\text{Cov}_{q_\theta(\omega)}[\mathbf{y}_\nu | \mathbf{x}_\nu])$  [22]. Additional measures of uncertainty will be presented in the next section.

*Step 3: Detecting adversarial inputs.* Given  $U(\mathbf{x}_\nu)$ , detection of adversarial inputs is cast as testing the hypotheses

$$\begin{cases} \mathcal{H}_0 : \mathbf{x}_\nu = \mathbf{x}_\nu^{\text{clean}} & U(\mathbf{x}_\nu) \leq \tau_0 \\ \mathcal{H}_1 : \mathbf{x}_\nu = \mathbf{x}_\nu^{\text{clean}} + \mathbf{n}_\nu^{\text{adv}} & U(\mathbf{x}_\nu) > \tau_0 \end{cases} \quad (8)$$

1. Such as back propagation based on e.g., a cross-entropy criterion.

where the null suggests absence of adversarial perturbation (low variance/uncertainty below threshold  $\tau_0$ ), while the alternative in effect raises a red flag for presence of adversarial input (high variance/uncertainty above threshold  $\tau_0$ ).

We will now proceed to introduce our novel variational distribution model targeting improved detection of adversaries based on uncertainty minimization.

### 3 MINIMUM UNCERTAINTY BASED DETECTION

To design  $q_z(\mathbf{z}_l; \boldsymbol{\theta}_{z_l})$ , we will build on and formalize the sampling scheme in [44] that is employed to specify the joint pmf of the (generally correlated) binary variables  $\{z_{l,i}\}_{i=1}^{h_l}$  per layer  $l$ . To this end, we randomly pick one activation unit output of the  $h_l$  hidden units per layer  $l$ ; and repeat such a random draw  $C$  times with replacement. Let  $\zeta_l^{(c)}$  denote per draw  $c$  the  $h_l \times 1$  vector variable

$$\zeta_l^{(c)} = [\zeta_{l,1}^{(c)}, \zeta_{l,2}^{(c)} \dots \zeta_{l,h_l}^{(c)}]^\top \sim \text{Categorical}(\mathbf{p}_l), \quad c = 1, \dots, C$$

where each entry  $\zeta_{l,i}^{(c)}$  is a binary random variable with

$$\zeta_{l,i}^{(c)} = \begin{cases} 1 & \text{if draw } c \text{ picks the } i\text{th unit of hidden layer } l \\ 0 & \text{otherwise} \end{cases}$$

and the  $h_l \times 1$  vector  $\mathbf{p}_l$  with nonnegative entries summing up to 1 specifies the Categorical pmf of  $\zeta_l^{(c)}$ .

With  $\|\cdot\|$  denoting element-wise binary OR operation on vectors  $\{\zeta_l^{(c)}\}_{c=1}^C$ , we define next the vector

$$\mathbf{z}_l := \zeta_l^{(1)} \|\zeta_l^{(2)}\| \dots \|\zeta_l^{(C)}\|. \quad (9)$$

Using  $\mathbf{z}_l$  as in (9) with  $\{\boldsymbol{\theta}_{z_l} = \mathbf{p}_l\}_{l=1}^L$  to be selected, enables finding the expectation and then  $\mathbf{M}_l$  in (6). Deterministic matrix  $\mathbf{M}_l$  along with the variates  $\{\mathbf{z}_l^{(r)}\}_{r=1}^R$  provide the desired DNN realizations to estimate the uncertainty  $U(\mathbf{x}_\nu; \{\mathbf{p}_l\}_{l=1}^L) = \text{Tr}(\text{Cov}_{q_\theta(\omega)}[\mathbf{y}_\nu | \mathbf{x}_\nu])$  as in (7). In turn, this leads to our novel adversarial input detector (cf. (8))

$$\begin{cases} \mathcal{H}_0 : \mathbf{x}_\nu = \mathbf{x}_\nu^{\text{clean}} & \min_{\{\mathbf{p}_l\}_{l=1}^L} U(\mathbf{x}_\nu; \{\mathbf{p}_l\}_{l=1}^L) \leq \tau_0 \\ \mathcal{H}_1 : \mathbf{x}_\nu = \mathbf{x}_\nu^{\text{clean}} + \mathbf{n}_\nu^{\text{adv}} & \text{otherwise} \end{cases} \quad (10)$$

where variational parameters  $\{\mathbf{p}_l\}_{l=1}^L$  are sought such that uncertainty  $U(\mathbf{x}_\nu; \{\mathbf{p}_l\}_{l=1}^L)$  is minimized under  $\mathcal{H}_0$ .

The rationale behind our detector in (10) is that for a given detection threshold  $\tau_0$ , uncertainty minimization will increase the number of clean images whose minimized uncertainty will fall below this threshold, and thus lead to a lower probability of false alarms. The probability of adversarial input detection however, depends on test statistic pdf under  $\mathcal{H}_1$ , in which the adversarial perturbation  $\mathbf{n}_\nu^{\text{adv}}$  is unknown. The premise here is that due to *network instability* under  $\mathcal{H}_1$ , the sought probabilities  $\{\mathbf{p}_l\}_{l=1}^L$  will not reduce uncertainty under  $\mathcal{H}_1$  as effectively, thus minimum uncertainty-based detection can provide improved ROC curves. In lieu of analytical metrics, this has been tested through extensive numerical experiments, and its effectiveness has been empirically corroborated.

Furthermore, Table 1 provides a list of variables and their definition to improve readability.

Param.	Definition
$\mathbf{x}_\nu$	Test input image
$\mathbf{y}_\nu^{\text{target}}$	Ground-truth one-hot class label for input $\mathbf{x}_\nu$
$\mathbf{y}_\nu$	Output for $\mathbf{x}_\nu$ by the deterministic classifier
$\mathbf{y}_\nu^{(r)}$	Output for $\mathbf{x}_\nu$ in the $r$ -th realization of the detection network (with sampling units)
$\mathbf{x}_{(l)}$	Value at hidden layer $l$ with width $h_l$
$\mathbf{W}_l$	Network weight at layer $l$ (random variable)
$\mathbf{W}_l^{\text{TR}}$	Network weights given by the training phase
$\mathbf{M}_l$	Expected value of $\mathbf{W}_l$
$\mathbf{S}_l$	Sampling matrix, defined as $\mathbf{S}_l = \text{diag}(\mathbf{z}_l)$
$\mathbf{D}_l$	Diagonal matrix $\mathbf{D}_l = \text{diag}^\dagger(\mathbb{E}_{q_z(\mathbf{z}_l; \mathbf{p}_l)}[\mathbf{z}_l])$
$\mathbf{z}_l$	$h_l$ -dimensional binary multivariate, modeling the overall sampling outcome at layer $l$
$\zeta_l^c$	Categorical $h_l$ -dimensional binary multivariate, modeling the $c$ -th draw at layer $l$
$\mathbf{p}_l$	Parameters of the categorical pmf of $\zeta_l^c$
$C$	Total number of draws
$f$	Scaler coeff. in $[0, 1]$ defining $C = f \times \text{nnz}(\mathbf{x}_l)$
$B$	Resnet architecture blocks in Tables 3 and 4

TABLE 1: List of variables

### 3.1 Uncertainty measures

In order to carry the hypothesis test in (10), one has options for  $U(\mathbf{x}_\nu; \{\mathbf{p}_l\}_{l=1}^L)$  other than the conditional variance. For DNNs designed for classification, mutual information has been recently proposed as a measure of uncertainty [38]

$$\hat{I}(\mathbf{x}_\nu; \{\mathbf{p}_l\}_{l=1}^L) := H(\bar{\mathbf{y}}_\nu) - \frac{1}{R} \sum_{r=1}^R H(\mathbf{y}_\nu^{(r)}) \quad (11)$$

where superscript  $r$  indexes the pass of input  $\mathbf{x}_\nu$  through the  $r$ th DNN realization with corresponding random output  $\mathbf{y}_\nu^{(r)} := [y_{\nu,1}^{(r)}, y_{\nu,2}^{(r)}, \dots, y_{\nu,K}^{(r)}]^\top$  in a  $K$ -class classification task, and  $H(\cdot)$  is the entropy function<sup>2</sup>

$$H(\mathbf{y}_\nu) := - \sum_{k=1}^K y_{\nu,k} \log(y_{\nu,k}). \quad (12)$$

The test statistic in (10) requires finding  $\{\mathbf{p}_l\}_{l=1}^L$  by solving

$$\min_{\{\mathbf{p}_l\}_{l=1}^L} \hat{I}(\mathbf{x}_\nu; \{\mathbf{p}_l\}_{l=1}^L) \quad (13)$$

which is highly non-convex. However, using Taylor's expansion of the logarithmic terms in (12), one can approximate the mutual information in (11) with the variance score  $\text{Tr}(\text{Cov}_{q_\theta(\omega)}[\mathbf{y}_\nu])$  in (10), where the conditioning on  $\mathbf{x}_\nu$  has been dropped for brevity [38]. As a result, the optimization in (13) is approximated as

$$\min_{\{\mathbf{p}_l\}_{l=1}^L} U(\mathbf{x}_\nu; \{\mathbf{p}_l\}_{l=1}^L) = \text{Tr}(\text{Cov}_{q_\theta(\omega)}[\mathbf{y}_\nu]). \quad (14)$$

To solve (14), one needs to express the objective in terms of the optimization variables  $\{\mathbf{p}_l\}$  for all layers explicitly. To this end, the following section studies a two-layer network, whose result will then be generalized to deeper models.

2. Entropy functions in (11) are also parameterized by  $\{\mathbf{p}_l\}_{l=1}^L$ , but we abbreviate them here as  $H(\bar{\mathbf{y}}_\nu)$  and  $H(\mathbf{y}_\nu^{(r)})$ .

### 3.2 Simplification of the predictive variance

Aiming at a convenient expression for the cost in (14), consider first a two-layer network with input-output (I/O) relationship<sup>3</sup>

$$\mathbf{y}_\nu = \sigma_{\text{softmax}}(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}_\nu)) \quad (15)$$

where  $\mathbf{W}_1, \mathbf{W}_2$  are random matrices corresponding to the weights of the two layers as in (6), while  $\sigma_{\text{softmax}}$  is the softmax memoryless nonlinearity

$$\sigma_{\text{softmax}}(\mathbf{u}) := \left[ \frac{e^{u_1}}{\sum_{i=1}^K e^{u_i}}, \frac{e^{u_2}}{\sum_{i=1}^K e^{u_i}}, \dots, \frac{e^{u_K}}{\sum_{i=1}^K e^{u_i}} \right]^\top$$

with  $\mathbf{u} := [u_1, u_2, \dots, u_K]^\top$ , and the inner  $\sigma$  in (15) models a general differentiable nonlinearity such as tanh. Although differentiability of the nonlinearities is needed for the derivations in this section, the general idea will be later tested on networks with non-differentiable nonlinearities (such as ReLU) in the experiments.

Given trained weights  $\{\mathbf{W}_l^{(TR)}\}_{l=1}^2$ , and using (4) and (6), the random weight matrices are found as

$$\mathbf{W}_l := \mathbf{M}_l \text{diag}(\mathbf{z}_l) = \mathbf{W}_l^{(TR)} \mathbf{S}_l \mathbf{D}_l \quad l = 1, 2 \quad (16)$$

where  $\mathbf{S}_l = \text{diag}(\mathbf{z}_l)$  denotes the random sampling matrix with pseudo-inverse diagonal mean given by  $\mathbf{D}_l = \text{diag}^\dagger(\mathbb{E}_{q_z(\mathbf{z}_l; \mathbf{p}_l)}[\mathbf{z}_l])$ . Since  $\mathbb{E}[\mathbf{W}_l] \mathbf{x}_{(l-1)} = \mathbf{W}_l^{(TR)} \mathbf{x}_{(l-1)}$ , the mean of  $\mathbf{W}_l$  does not depend on  $\mathbf{p}_l$ , while its higher-order moments do.

**Proposition 1.** For the two-layer network in (15), the proposed minimization in (14) can be approximated by

$$\min_{\{\mathbf{p}_i \geq \mathbf{0}, \mathbf{1}^\top \mathbf{p}_i = 1\}_{i=1}^2} \text{Tr}(\text{Cov}_{\mathbf{W}_2}[\mathbf{W}_2 \sigma(\mathbf{W}_1^{(TR)} \mathbf{x}_\nu)]) \quad (17)$$

$$+ \gamma \text{Tr}(\mathbb{E}_{\mathbf{W}_2}[\mathbf{W}_2 \mathbf{W}_2^\top]) \text{Tr}(\text{Cov}_{\mathbf{W}_1}[\mathbf{W}_1 \mathbf{x}_\nu])$$

where  $\gamma$  is a constant. The solution of (17) proceeds in two steps

**Step 1:**  $\mathbf{p}_1^* = \arg \min_{\mathbf{p}_1} \text{Tr}(\text{Cov}_{\mathbf{W}_1}[\mathbf{W}_1 \mathbf{x}_\nu])$

**Step 2:**  $\mathbf{p}_2^* = \arg \min_{\mathbf{p}_2} \text{Tr}(\text{Cov}_{\mathbf{W}_2}[\mathbf{W}_2 \sigma(\mathbf{W}_1^{(TR)} \mathbf{x}_\nu)])$   
 $+ \gamma' \text{Tr}(\mathbb{E}_{\mathbf{W}_2}[\mathbf{W}_2 \mathbf{W}_2^\top])$

where  $\gamma' := \gamma \text{Tr}(\text{Cov}_{\mathbf{W}_1}[\mathbf{W}_1 \mathbf{x}_\nu]) \Big|_{\mathbf{p}_1 = \mathbf{p}_1^*}$ .

*Proof.* See Appendix 8.1.

**Remark.** The cost in (17) approximates that in (14) by casting the overall uncertainty minimization as a weighted sum of layer-wise variances. In particular,  $\mathbf{p}_1^*$  is the sampling probability vector that minimizes variance score of the first layer. It subsequently influences the regularization scalar  $\gamma'$  in minimizing the second layer variance, which yields the pmf vector  $\mathbf{p}_2^*$ . This can be inductively generalized to  $L > 2$  layers. As  $L$  increases however, so do the number of cross terms. For simplicity and scalability, we will further approximate the per-layer minimization by dropping the regularization term, which leads to *separable* optimization

3. Derivations in this section carry over readily to a more general I/O  $\mathbf{y}_\nu = \sigma_{\text{softmax}}(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}_\nu + \mathbf{b}_1) + \mathbf{b}_2)$  with  $\mathbf{b}_1$  and  $\mathbf{b}_2$  deterministic.

across layers. This is an intuitively pleasing relaxation, because layer-wise variance is minimized under  $\mathcal{H}_0$ , which also minimizes the regularization weight  $\gamma'$ .

The resulting non-regularized approximant of step 2 is

$$\mathbf{p}_2^* = \arg \min_{\mathbf{p}_2} \text{Tr}(\text{Cov}_{\mathbf{W}_2}[\mathbf{W}_2 \sigma(\mathbf{W}_1^{(TR)} \mathbf{x}_\nu)])$$

generalizing to the  $l$ -th layer in an  $L$ -layer DNN as

$$\mathbf{p}_l^* = \arg \min_{\mathbf{p}_l} \text{Tr}(\text{Cov}_{\mathbf{W}_l}[\mathbf{W}_l \mathbf{x}_{(l-1)}]) \quad (18)$$

where  $\mathbf{x}_{(l-1)}$  is the output of the  $(l-1)$ st layer, regardless of pmf vectors of other layers  $\{\mathbf{p}_{l'}\}_{l' \neq l}$ .

### 3.3 Layer-wise variance minimization

Here we will solve the layer-wise variance minimization in (18). Using (16), the cost can be upper bounded by

$$\begin{aligned} & \text{Tr}(\text{Cov}_{\mathbf{W}_l}[\mathbf{W}_l \mathbf{x}_{(l-1)}]) \\ &= \mathbb{E} \left[ \left\| \mathbf{W}_l^{(TR)} \mathbf{S}_l \mathbf{D}_l \mathbf{x}_{(l-1)} - \mathbb{E}[\mathbf{W}_l^{(TR)} \mathbf{S}_l \mathbf{D}_l \mathbf{x}_{(l-1)}] \right\|_2^2 \right] \\ &= \mathbb{E} \left[ \left\| \mathbf{W}_l^{(TR)} \mathbf{S}_l \mathbf{D}_l \mathbf{x}_{(l-1)} - \mathbf{W}_l^{(TR)} \mathbf{x}_{(l-1)} \right\|_2^2 \right] \\ &\leq \left\| \mathbf{W}_l^{(TR)} \right\|_2^2 \mathbb{E} \left[ \left\| \mathbf{S}_l \mathbf{D}_l \mathbf{x}_{(l-1)} - \mathbf{x}_{(l-1)} \right\|_2^2 \right] \\ &= \left\| \mathbf{W}_l^{(TR)} \right\|_2^2 \sum_{i=1}^{h_l} \mathbb{E}[(S_{ii} D_{ii} \mathbf{x}_{(l-1),i} - \mathbf{x}_{(l-1),i})^2] \\ &= \left\| \mathbf{W}_l^{(TR)} \right\|_2^2 \sum_{i=1}^{h_l} \mathbf{x}_{(l-1),i}^2 \mathbb{E}[(S_{ii} D_{ii} - 1)^2] \\ &= \left\| \mathbf{W}_l^{(TR)} \right\|_2^2 \sum_{i=1}^{h_l} \mathbf{x}_{(l-1),i}^2 \left( \frac{1}{\pi_{l,i}} - 1 \right) \end{aligned} \quad (19)$$

where the last equality follows because the  $C$  draws are iid with replacement, and the binary random variables  $z_{l,i}$  reduce to Bernoulli ones with parameter  $\pi_{l,i} = 1 - (1 - p_{l,i})^C$ ; hence, for  $\mathbf{x}_{(l-1),i} \neq 0$  it holds that  $\mathbb{E}[S_{ii}^2 D_{ii}^2] = 1/\pi_{l,i}$  and  $\mathbb{E}[S_{ii} D_{ii}] = 1$ , which implies that  $\mathbb{E}[(S_{ii} D_{ii} - 1)^2] = (1/\pi_{l,i}) - 2 + 1$ .

Using (19), the optimization in (18) can be approximately solved by a majorized surrogate as

$$\min_{\mathbf{p} \geq \mathbf{0}, \mathbf{1}^\top \mathbf{p} = 1} \sum_{i=1}^{h_l} \frac{1}{1 - (1 - p_{l,i})^C} \mathbf{x}_{(l-1),i}^2 \quad (20)$$

which is a convex problem that can be solved efficiently as elaborated next.

## 4 SOLVING LAYER-BY-LAYER MINIMIZATION

Consider rewriting the layer-wise variance minimization in (20) in a general form as

$$\min_{\mathbf{p} \geq \mathbf{0}, \mathbf{1}^\top \mathbf{p} = 1} \sum_{i=1}^h \frac{\alpha_i}{1 - (1 - p_i)^C} \quad (21)$$

where  $\alpha_i := \mathbf{x}_{(l-1),i}^2$  for the  $l$ -th layer. Over the feasible set of the probability simplex, the cost in (21) has semi-definite Hessian; thus, it is convex, and can be solved by projected gradient descent iterations. However,  $\mathbf{p}$  lies in the probability simplex space of dimension  $h$ , the number of hidden nodes in a given layer, and is typically very large.

The large number of variables together with possible ill-conditioning can slow down the convergence rate.

To obtain a solver with quadratic convergence rate, we build on the fact that  $h_l$  is usually very large, which implies that  $p_i \ll 1$  for the practical setting at hand. Using the inequality  $1 - (1 - p_i)^C \geq 1 - e^{-Cp_i}$ , the cost in (21) can then be tightly upperbounded, which leads to majorizing (21) as

$$\min_{\mathbf{p} \geq \mathbf{0}, \mathbf{1}^\top \mathbf{p} = 1} \sum_{i=1}^h \frac{\alpha_i}{1 - e^{-Cp_i}} \quad (22)$$

The KKT conditions yield the optimal solution of the convex problem in (22), as summarized next.

**Proposition 2.** *The optimization in (22) can be solved with quadratic convergence rate, and the optimum is given by*

$$p_i^* = -\frac{1}{C} \ln \left( \frac{2\rho^* + x_{(l-1),i}^2 - \sqrt{[2\rho^* + x_{(l-1),i}^2]^2 - 4\rho^{*2}}}{2\rho^*} \right) \quad (23)$$

where  $\rho^*$  is the solution to the following root-finding problem

$$\sum_{i=1}^h \ln(2\rho + x_{(l-1),i}^2 - \sqrt{[2\rho + x_{(l-1),i}^2]^2 - 4\rho^2}) - n \ln(2\rho) + C = 0.$$

*Proof.* See Appendix 8.2.

### 4.1 Approximate variance minimization for small $C$

For small values of  $C$ , it holds that  $(1 - p_i)^C > 1 - Cp_i$ ; hence, the Bernoulli parameter  $\pi_i = 1 - (1 - p_i)^C$  can be approximated by its upperbound  $Cp_i > \pi_i$ . With this we can approximate the cost in (20), as

$$\min_{\mathbf{p} \geq \mathbf{0}, \mathbf{1}^\top \mathbf{p} = 1} \sum_{i=1}^h \frac{\alpha_i}{Cp_i} \quad (24)$$

Using the Lagrangian and the KKT conditions, we then find  $p_i^* = \sqrt{\alpha_i} / \sum_j \sqrt{\alpha_j}$ , which for the  $l$ -th layer is expressible as

$$p_{(l-1),i}^* = \frac{|x_{(l-1),i}|}{\sum_{j=1}^{h_l} |x_{(l-1),j}|} \quad (25)$$

This approximation provides analytical justification for the heuristic approach in [44], where it is proposed to sample with probabilities proportional to the magnitude of the hidden unit outputs. However, there remains a subtle difference, which will be clarified in Section 6.

Approximating (22) with (24) can be loose for large values of  $C$ , which motivates our next approximation.

### 4.2 Approximate variance minimization for large $C$

Building on the tight approximation in (22), one can further approximate the variance for large  $C$  as

$$\min_{\mathbf{p} \geq \mathbf{0}, \mathbf{1}^\top \mathbf{p} = 1} \sum_{i=1}^h \frac{\alpha_i}{1 - e^{-Cp_i}} \simeq \min_{\mathbf{p} \geq \mathbf{0}, \mathbf{1}^\top \mathbf{p} = 1} \sum_{i=1}^h \alpha_i (1 + e^{-Cp_i})$$

where we have used  $(1 - \delta)^{-1} \simeq 1 + \delta$  as a tight approximation for  $0 < \delta \ll 1$ . This leads to the minimization

$$\min_{\mathbf{p} \geq \mathbf{0}, \mathbf{1}^\top \mathbf{p} = 1} \sum_{i=1}^n \alpha_i e^{-Cp_i}$$

which again is a convex problem, whose solution can be obtained using the KKT conditions that lead to

$$-C\alpha_i e^{-C\hat{p}_i^*} + \lambda = 0 \quad \forall i$$

where  $\lambda$  is the Lagrange multiplier. Under the simplex constraint on the  $\{p_i\}$ , this leads to the optimal

$$\hat{p}_i^* = \left[ \frac{1}{C} \ln [Cx_{(l-1),i}^2] + \hat{\beta}^* \right]_+ \quad (26)$$

with  $[\cdot]_+$  denoting the projection on the positive orthant, and the normalization constant  $\beta := -\ln \lambda / C$  having optimal value

$$\hat{\beta}^* = \frac{C - \sum_{i=1}^{h_l} \ln(Cx_{(l-1),i}^2) \mathbf{1}_{\{\hat{p}_i^* > 0\}}}{C \sum_{i=1}^{h_l} \mathbf{1}_{\{\hat{p}_i^* > 0\}}}.$$

Although the solution to the fixed point condition cannot be obtained at one shot, and may require a few iterations to converge, in practice we only perform it once and settle with the obtained approximate solution  $\{\hat{p}_i^*\}_{i=1}^{h_l}$ .

## 5 PRACTICAL ISSUES

The present section deals with efficient implementation of the proposed approach in practice, and establishes links with state-of-the-art randomization-based detection methods.

### 5.1 Efficient implementation via non-uniform dropout

The proposed defense builds on modeling the variational pdf  $q_\theta(\omega)$  using a sampling-with-replacement process. Performing the proposed process however, may incur overhead complexity during inference when compared to the inexpensive dropout alternative outlined in Sec. 2.1. To reduce this complexity, one can implement our approach using efficient approximations, while leveraging the sampling probabilities learned through our uncertainty minimization.

Reflecting on the binary variables  $\{z_{l,i}\}$  that model the pickup of the hidden node  $i$  in the overall sampling process in (9), one can approximate the joint pmf of  $\{z_{l,i}\}_{i=1}^{h_l}$  as

$$q_z(\mathbf{z}_l; \mathbf{p}_l) \simeq \prod_{i=1}^{h_l} q_z(z_{l,i}; p_{l,i}) \quad (27)$$

where random variables  $\{z_{l,i}\}_i$  are now viewed as approximately independent *non-identical* Bernoulli variables with parameters  $\{\pi_{l,i}\}_{i=1}^{h_l}$ ; that is,  $z_{l,i} \sim \text{Bernoulli}(\pi_{l,i})$  for  $i = 1, \dots, h_l$ , where  $\pi_{l,i} = 1 - (1 - p_{l,i})^C$ .

Although (27) is an approximation, it provides insight but also an efficient implementation of the sampling process. In fact, the proposed optimization in (21) can now be viewed as an optimization over the non-uniform dropout probabilities, coupled implicitly through the hyper-parameter  $C$ , whose selection guarantees a certain level of randomness. This is to be contrasted with finding optimal dropout probabilities - a task requiring grid search over an  $h_l$ -dimensional space for layer  $l$ , where  $h_l$  can be hundreds of thousands to millions in CNNs classifying high-quality images. Interestingly, the proposed convex optimization simplifies the high-dimensional grid-search into a scalar root-finding task, whose solution can be efficiently found with super-linear (quadratic) convergence rate.

### 5.2 Placement and adjustment of the sampling units

It has been argued that CNN layers at different depths can provide extracted features with variable levels of expressiveness [48]. On a par with this, one can envision the defense potential at different depths by incorporating sampling units across say  $B$  blocks of the network as listed in Tables 3 and 4. In particular, the dropout defense has been mostly utilized at the last layer after flattening [38], whereas here we consider the potential of sampling at earlier layers that has gone mostly under-explored so far. This can in turn result in DNN-based classifiers with robustness to adversarial attacks, as optimal sampling at the initial layers maybe crucial for correct detection of the adversarial input. We henceforth refer to a DNN (or CNN) equipped with random sampling as the *detection network*, and the original one without the sampling units as the *full network*.

Similar to the pick up probability  $\pi$  in uniform dropouts, the number of draws  $C$  in our approach is a hyper parameter that controls the level of randomness present in the detection network. Qualitatively speaking, the smaller number of units (smaller  $C$ ) is picked per layer, the larger ‘amount of randomness’ emerges (further  $\pi_{l,i}$  is from 1). This can lead to forward propagating not as informative (under-sampled) features, meaning not representative of the clean image, and can thus cause unreliable detection. A large  $C$  on the other hand, increases the probability to pick up units per layer, which requires a large number of MC realizations for reliable detection, otherwise small randomness will lead to miss-detection. At the extreme, very large  $C$  renders the detection and full networks identical, thus leading to unsuccessful detection of adversarial inputs. In a nutshell, there is a trade-off in selecting  $C$ , potentially different for the initial, middle, and final layers of a given CNN.

Fig. 1 categorizes existing and the novel randomization-based approaches to detecting adversarial inputs.

**Uniform dropout.** In this method, units are independently dropped w.p.  $1 - \pi$ , and sampled (picked) w.p.  $\pi \forall l, i$ .

**Non-uniform dropout using variance minimization.** Dropout here follows the scheme in subsection 5.1, for which we pursue the following two general cases with *deterministic* and *dynamic* probabilities.

(C1) *Variance minimization with fixed probabilities.* In this case, the image is first passed through the full network to obtain the values  $\{x_{(l-1),i}\}$  of the unit outputs per hidden layer. These are needed to determine the non-uniform dropout probabilities  $1 - p_{l,i}$  (thus  $\pi_{l,i}$  and then the index of the units to sample) via *exact*, *linear*, or *logarithmic* approximations given in (23), (25) and (26), respectively, referred to as VM-exact, VM-lin, and VM-log; see Fig. 2-a.

Despite parallel MC passes in the proposed class of sampling with fixed probabilities (step 3 in Fig. 2-a), the first step still imposes a *serial* overhead in detection since the wanted probabilities must be obtained using a pass through the full network. Our approach to circumventing this overhead is through approximation using the following class of sampling with *dynamic* probabilities.

(C2) *Variance minimization with dynamic probabilities.* Rather than finding the sampling probabilities beforehand,  $p_{l,i}^{(r)}$  are determined *on-the-fly* as the image is passed through the detection network with the units sampled per layer.

## Randomization-based approaches to detecting adversaries in DNNs

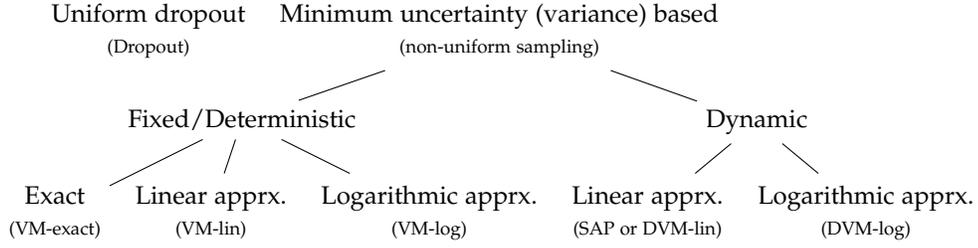
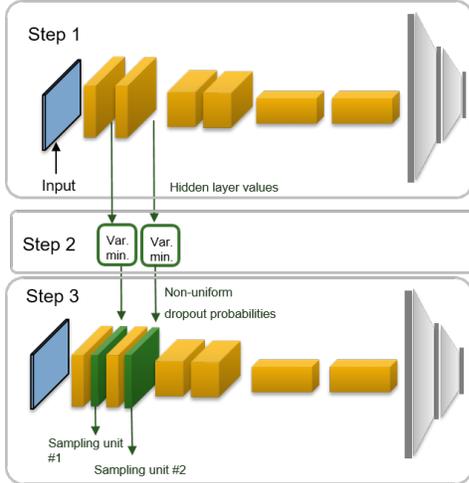
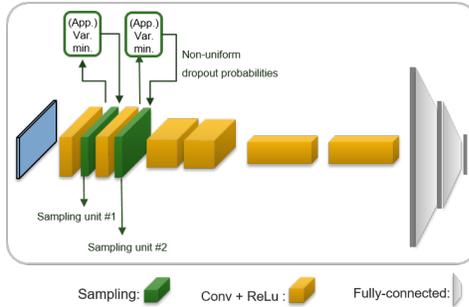


Fig. 1. Overview of randomization-based adversary detection schemes



a) Detection via deterministic sampling probabilities



b) Detection via dynamic sampling probabilities

Fig. 2. Schematic of the proposed detection schemes

As a result, the observed unit values are random (after passing through at least one unit sampled), and are different across realizations. In order to mitigate solving many optimization problems, variance minimization with dynamic probabilities is only implemented via *linear* and *logarithmic* approximations (25) and (26); and are referred to as DVM-lin and DVM-log, respectively; see Fig. 2-b.

It is interesting to note that DVM-lin corresponds to the proposed stochastic activation pruning (SAP) in [44], with

$$\mathbf{p}_{l,r}^{\text{SAP}} = \left[ \frac{|x_{(l-1),1}^{(r)}|}{\sum_{i=1}^{h_l} |x_{(l-1),i}^{(r)}|}, \frac{|x_{(l-1),2}^{(r)}|}{\sum_{i=1}^{h_l} |x_{(l-1),i}^{(r)}|}, \dots, \frac{|x_{(l-1),h_l}^{(r)}|}{\sum_{i=1}^{h_l} |x_{(l-1),i}^{(r)}|} \right]$$

where  $x_{(l-1),i}^{(r)}$  is the output of the  $i$ -th activation unit of the

**Algorithm 1:** Adversary detection - fixed  $\{p_{l,i}\}$ 

**Input:** Test image  $\mathbf{x}_\nu$ ,  $B, C, R$  and  $\tau_0$

- 1 Pass image  $\mathbf{x}_\nu$  through full network; find  $\{x_{(l-1),i}\}$
  - 2 Use  $\{x_{(l-1),i}\}$  to obtain  $\{p_{l,i}\}$  via (23), (25) or (26)
  - 3 **for**  $r = 1, 2, \dots, R$  **do**
  - 4     Collect output class  $\mathbf{y}_\nu^{(r)}$
  - 5 **end**
  - 6 Estimate the mutual information (MI) of  $\{\mathbf{y}_\nu^{(r)}\}_{r=1}^R$
- Output:** Declare *adversary* if MI exceeds threshold  $\tau_0$

**Algorithm 2:** Adversary detection - dynamic  $\{p_{l,i}\}$ 

**Input:** Test image  $\mathbf{x}_\nu$ ,  $B, C, R$  and  $\tau_0$

- 1 **for**  $r = 1, 2, \dots, R$  **do**
  - 2     Collect  $\mathbf{y}_\nu^{(r)}$  after passing  $\mathbf{x}_\nu$  through the detection network with units picked with dynamic probabilities obtained (exactly or approximately) using the observed values
  - 3 **end**
  - 4 Estimate the mutual information (MI) of  $\{\mathbf{y}_\nu^{(r)}\}_{r=1}^R$
- Output:** Declare *adversary* if MI exceeds threshold  $\tau_0$

$l$ -th layer in the  $r$ -th realization for input  $\mathbf{x}$ .

Figure 1 provides an overview of the sampling methods, while Algorithms 1 and 2 outline the two proposed variance minimization-based detection methods in pseudocode.

## 6 NUMERICAL TESTS

**Algorithm 3:** Layer-wise minimum variance solver

- 1 **Solve:**  $\min_{\mathbf{p} \geq 0, \mathbf{1}^\top \mathbf{p} = 1} \sum_{i=1}^h \frac{\alpha_i}{1 - (1 - p_i)^C}$
- Input :**  $[\alpha_1, \alpha_2, \dots, \alpha_h], C$
- Output:** Nonuniform dropout pmf  $\boldsymbol{\pi} = [\pi_1 \dots \pi_h]^\top$
- 2 Using bisection and initialization  $\rho_0 = \sum_{i=1}^h \alpha_i/h$ , find the root  $\rho^*$  for 
$$\sum_i \ln(2\rho + \alpha_i - \sqrt{(2\rho + \alpha_i)^2 - 4\rho^2}) - n \ln(2\rho) + C = 0$$
- 3 Set  $p_i^* = -\frac{1}{C} \ln\left(\frac{2\rho^* + \alpha_i - \sqrt{(2\rho^* + \alpha_i)^2 - 4\rho^{*2}}}{2\rho^*}\right) \quad \forall i$
- 4 Set  $\pi_i^* = 1 - (1 - p_i^*)^C \quad \forall i$

Dataset	image size	# train	# val.	# test
CIFAR10	32 x 32	50,000	2,000	8,000
Cats-and-dogs	224 x 224	10,000	2,000	13,000

TABLE 2: CIFAR10 and cats-and-dogs image-classification datasets

name	output-size	20 layers	#sampling units
Block1	32 x 32	[3 x 3, 16]	1
Block2	32 x 32	$\begin{matrix} 3 \times 3, 16 \\ 3 \times 3, 16 \end{matrix} \times 3$	6
Block3	16 x 16	$\begin{matrix} 3 \times 3, 32 \\ 3 \times 3, 32 \end{matrix} \times 3$	6
Block4	8 x 8	$\begin{matrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{matrix} \times 3$	6
Block5	1 x 1	average pool, 64-d fully conn., softmax	1

TABLE 3: ResNet20 architecture on CIFAR10 dataset

In this section, we test the effectiveness of the proposed sampling method for detecting various adversarial attacks on CNNs used for image classification. In order to address the raised issue in [13], classification of the CIFAR10 image dataset using ResNet20 as well as the high-resolution cats-and-dogs images using ResNet34 networks [49] are tested. A short summary of the two networks and datasets can be found in Tables 2, 3 and 4. In order to investigate the issue around placement of the sampling units, we will place them after ReLU activation layers in different “blocks” ( $B$ ) of the ResNet20 and ResNet34 networks, as listed in Tables 3 and 4. Numerical tests are made available online.<sup>4</sup>

### 6.1 CIFAR10 dataset

ResNet20 is trained using 20 epochs with minibatches of size 128. Adversarial inputs are crafted on the corresponding MC network as in [38], using the fast gradient

4. <https://github.com/FatemehSheikholeslami/variance-minimization>

	output-size	34 layers	#sampling units
Block 1	112 x 112	[7 x 7, 64], 3x3 max-pool	2
Block2	56 x 56	$\begin{matrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{matrix} \times 3$	6
Block3	28 x 28	$\begin{matrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{matrix} \times 4$	8
Block4	14 x 14	$\begin{matrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{matrix} \times 6$	12
Block5	7 x 7	$\begin{matrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{matrix} \times 3$	6
Block 6	1 x 1	average pool, 1000-d fc, softmax	1

TABLE 4: ResNet34 architecture on cats-and-dogs dataset

sign method (FGSM) [50], the basic iterative method (BIM) [51], the momentum iterative method (MIM) [52], and the Carlini-and-Wagner (C&W) [14] attacks. Parameters of the attacks as well as test accuracy of the MC network on clean and adversarial inputs are reported in Table 5.

Placement parameter  $B$  and sampling parameters  $C$  for variance minimization methods as well as the dropout probability for uniform dropout are selected by cross validation. To clarify the suboptimality gap between the exact and approximate variance minimization with **deterministic** sampling probabilities, we have cross-validated the parameters for VM-exact, and reused them for VM-lin and VM-log approximates.

The sampling parameter is selected as  $C = f \times \text{nnz}(x_l)$  for the  $l$ -th layer sampling unit, where  $\text{nnz}(\cdot)$  denotes the number of non-zero entries<sup>5</sup>, and  $f$  is the sampling ratio varied in  $f \in \{0.6, 0.7, 0.8, 0.9, 1.0, 1.5, 2.0, 3.0, 4.0\}$ .<sup>6</sup> Probability in uniform dropout is also varied as  $\pi_{\text{drop}} \in \{0.1, 0.2, \dots, 0.7\}$ , and the number of MC runs is  $R = 20$ .

In order to properly evaluate accuracy in detection of adversarial images, we only aim at detecting the test samples that are correctly classified by the full network, and misclassified after the adversarial perturbation. The detection performance is then reported in terms of the receiver operating characteristic (ROC) curve in Fig. 3, obtained by varying the threshold parameter  $\tau_0$ . The exact area-under-curve values along with parameters  $B, f, \pi_{\text{drop}}$  are also reported in Tables 6 and 7, highlighting the improved detection via the proposed variance minimization approach.

Furthermore, in order to target more realistic scenarios, where attack generation is unknown and may indeed be crafted via various methods, we have also tested the performance against a “combination attack,” in which the adversarial input crafted with all 7 settings of attacks are considered. This indeed corroborates that placement of the sampling units in the fourth block along with careful tuning of the sampling probabilities via VM-exact provides the highest curve against combination of attacks, while its approximations follow in performance, outperforming uniform dropout. For further discussion on sensitivity against parameter selection, see Appendix 8.3.

### 6.2 Cats-and-dogs dataset

Tests are also carried out for the cats-and-dogs dataset,<sup>7</sup> which consists of high-quality images classified into binary classes of cats and dogs. Images are resized to  $224 \times 224$ , and are classified using ResNet34 [49]. Weights of the convolutional layers are transferred from the network trained on the ImageNet dataset.<sup>8</sup> This is subsequently followed by a dropout,  $1000 \times 2$  fully-connected and softmax layer, whose weights are trained using 10,000 images; see Table 2. The FGMS, BIM, MIM, an C&W attacks are crafted, and parameters are reported in Table 8. Detection parameters are similarly selected by using the validation set and varying

5. This selection is chosen by taking into account the fact that, only non-zero samples will be dropped upon not being selected, while zero entries will remain unchanged regardless of the sampling outcome.

6. Since the sampling procedure is modeled with replacement, fraction  $f$  may be selected greater than 100%.

7. <https://www.microsoft.com/en-us/download/details.aspx?id=54765>

8. [https://github.com/qubvel/classification\\_models](https://github.com/qubvel/classification_models)

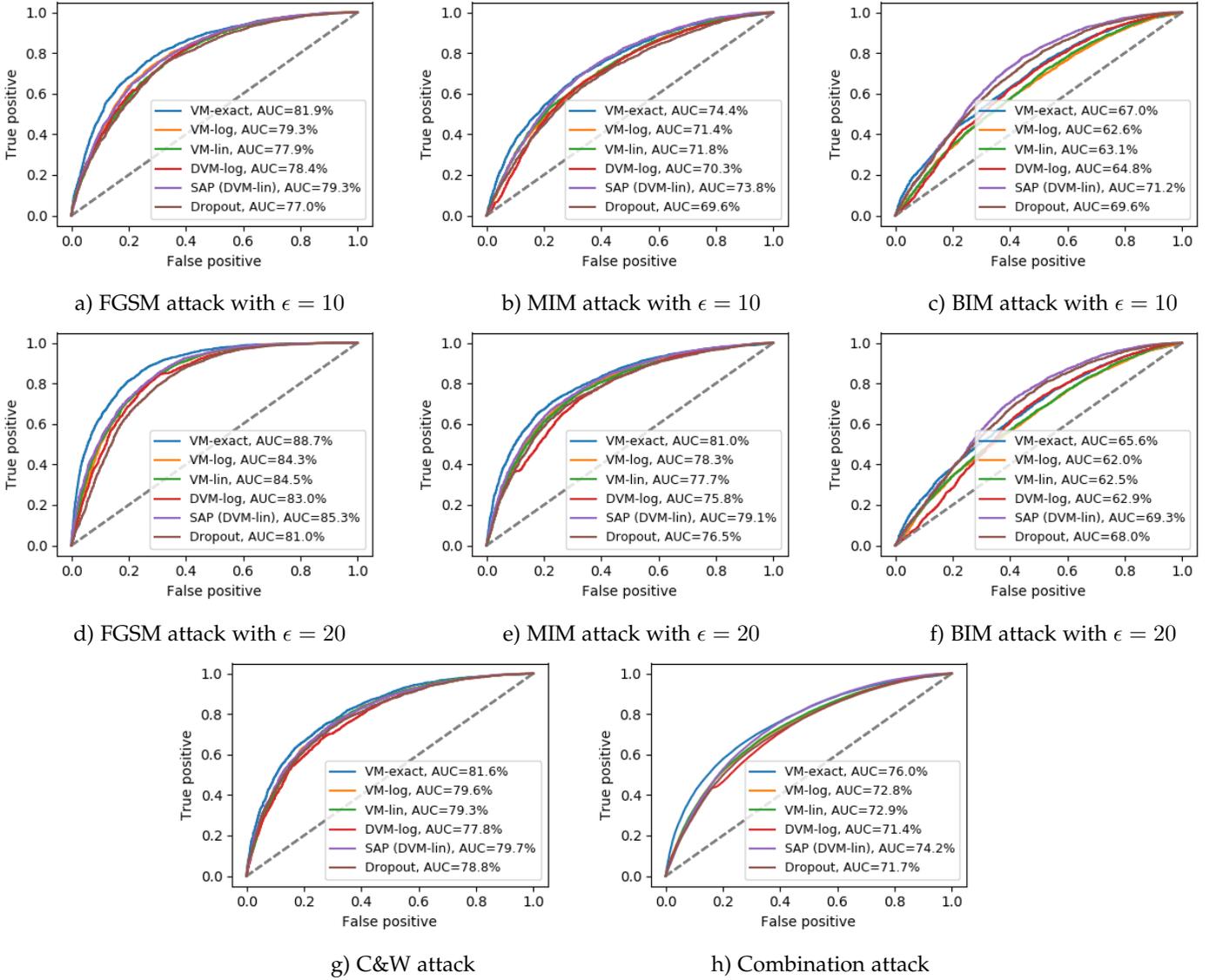


Fig. 3. ROC-curve of different attack-detection sampling schemes on CIFAR10 dataset against different attacks.

$B \in \{1, 2, 3, 4, 5, 6\}$ ,  $f \in \{0.6, 0.7, 0.8, 0.9, 1.0, 1.5, 2.0, 3.0\}$ , where  $C = f \times \text{nnz}(\mathbf{x}_l)$ ,  $\pi_{\text{drop}} \in \{0.1, 0.2, \dots, 0.7\}$ , and the number of MC runs is  $R = 20$ .

Fig. 4 plots the ROC curve for detection of adversarial versus clean images, and defense parameters are reported in Tables 9 and 7, quantifying the accuracy of attack detection across different methods. As with CIFAR10, tests are also extended to a combination attack, where detection is performed against the combination of all seven attacks with a fixed set of defense parameters  $B$ ,  $f$ , and  $\pi_{\text{drop}}$ .

Interestingly, it is observed that for small values of  $f$ , the linear approximation for variance minimization (VM-lin) follows the performance of the exact variance minimization (VM-exact) closely for FGSM, MIM, and C&W attacks, whereas the logarithmic approximation (VM-log) exhibits a large gap in performance. In contrast, for large values of  $f$ , VM-log demonstrates a smaller optimality gap with VM-exact as opposed to VM-lin; see Figs. 4 (c) and (f). This corroborates our approximations in (24) and (26), providing high-performance low-complexity substitutes for the exact

variance-minimization solver in both small and large sampling regimes, that is  $f < 1$  and  $f > 1$ . Similarly, improved performance of the logarithmic approximates versus the linear ones are also corroborated in the high-quality cats-and-dogs images versus CIFAR10, due to higher  $C$  resulting from higher dimensional vectors  $\mathbf{x}_{(l)}$  in the hidden layers.

The ROC curves further demonstrate that performance of the deterministic sampling probabilities, obtained by passing the image through the full network, are often superior to the dynamic ones (SAP and DVM-log), among which DVM-log demonstrates better performance.

### 6.3 Detection of adaptive attacks

In order to further evaluate the performance of the proposed detection schemes against white-box adaptive attacks, we now consider an adaptive attack setting, in which the attacker is aware of the defense mechanism, and designs the adversarial perturbation to jointly fool the classifier and the detector.

	clean	FGSM		BIM		MIM		C&W
Attack parameters	-	-		norm: $\infty$ # iter: 20 $\epsilon_{iter}: 1/255$		norm: $\infty$ # iter: 20 $\epsilon_{iter}: 1/255$		# binary search: 10 #max iter: 20 learning rate:0.1 initial const.: 10
		$\epsilon = 10/255$	$\epsilon = 20/255$	$\epsilon = 10/255$	$\epsilon = 20/255$	$\epsilon = 10/255$	$\epsilon = 20/255$	
Class. Acc.	91.5%	64.87%	56.91%	5.2%	5.0%	5.4%	5.1%	11.7%

TABLE 5: Attack parameters and test accuracy on clean and adversarial input in CIFAR10 dataset.

Sampling Method	FGSM Attack				MIM Attack			
	$\epsilon = 10$		$\epsilon = 20$		$\epsilon = 10$		$\epsilon = 20$	
	Parameters	AUC	Parameters	AUC	Parameters	AUC	Parameters	AUC
VM		<b>81.9</b>		<b>88.7</b>		<b>74.4</b>		<b>81.0</b>
VM-log	$(B, f) = (4, 2.0)$	79.3	$(B, f) = (4, 4.0)$	84.3	$(B, f) = (4, 4.0)$	71.4	$(B, f) = (4, 4.0)$	78.3
VM-linear		77.9		84.5		71.8		77.7
DVM-log	$(B, f) = (4, 3.0)$	78.4	$(B, f) = (5, 0.7)$	83.0	$(B, f) = (5, 4.0)$	70.3	$(B, f) = (4, 4.0)$	75.8
SAP	$(B, f) = (2, 4.0)$	79.3	$(B, f) = (3, 4.0)$	85.3	$(B, f) = (2, 3.0)$	73.8	$(B, f) = (3, 4.0)$	79.1
Dropout	$(B, \pi_{dnp}) = (5, 0.1)$	77.0	$(B, \pi_{dnp}) = (5, 0.1)$	81.0	$(B, \pi_{dnp}) = (5, 0.1)$	69.6	$(B, \pi_{dnp}) = (5, 0.2)$	76.5

TABLE 6: AUC-ROC of different attack-detection sampling schemes on CIFAR10 test set against FGSM and MIM attacks. Higher values indicate better detection.

Sampling Method	BIM Attack				C&W Attack		Combination Attack	
	$\epsilon = 10$		$\epsilon = 20$		parameters	AUC	parameters	AUC
	parameters	AUC	parameters	AUC				
VM		67.0		65.6		<b>81.6</b>		<b>76.0</b>
VM-log	$(B, f) = (4, 4.0)$	62.6	$(B, f) = (4, 4.0)$	62.0	$(B, f) = (4, 3.0)$	79.6	$(B, f) = (4, 4.0)$	72.8
VM-linear		63.1		62.5		79.3		72.9
DVM-log	$(B, f) = (1, 3.0)$	64.8	$(B, f) = (1, 4.0)$	62.9	$(B, f) = (5, 0.8)$	77.8	$(B, f) = (4, 3.0)$	71.4
SAP	$(B, f) = (2, 1.5)$	<b>71.2</b>	$(B, f) = (2, 1.5)$	<b>69.3</b>	$(B, f) = (5, 4.0)$	79.7	$(B, f) = (2, 3.0)$	74.2
Dropout	$(B, \pi_{dnp}) = (2, 0.1)$	69.6	$(B, \pi_{dnp}) = (2, 0.1)$	68.0	$(B, \pi_{dnp}) = (5, 0.2)$	78.8	$(B, \pi_{dnp}) = (5, 0.1)$	71.7

TABLE 7: AUC-ROC of different attack-detection sampling schemes on CIFAR10 test set against FGSM, C&W, and combination attacks. Higher values indicate better detection.

	clean	FGSM		BIM		MIM		C&W
Attack parameters	-	-		norm: $\infty$ # iter: 10 $\epsilon_{iter}: 1.5$		norm: $\infty$ # iter: 20 $\epsilon_{iter}: 1.5$		# binary search steps: 10 #max iter: 20 learning rate: 0.1 initial const.: 10
		$\epsilon = 10$	$\epsilon = 20$	$\epsilon = 10$	$\epsilon = 20$	$\epsilon = 10$	$\epsilon = 20$	
Class. Acc.	94.5%	74.85%	70.5%	19.2%	18.4%	12.9%	9.9%	68.95%

TABLE 8: Attack parameters and test accuracy on clean and adversarial input in cats-and-dogs dataset.

Sampling Method	FGSM Attack				MIM Attack			
	$\epsilon = 10$		$\epsilon = 20$		$\epsilon = 10$		$\epsilon = 20$	
	Parameters	AUC	Parameters	AUC	Parameters	AUC	Parameters	AUC
VM		73.5		84.0		<b>72.1</b>		78.5
VM-log	$(B, f) = (2, 0.7)$	57.2	$(B, f) = (2, 0.7)$	63.8	$(B, f) = (2, 0.7)$	58.1	$(B, f) = (2, 0.7)$	61.5
VM-linear		70.8		82.2		68.8		77.1
DVM-log	$(B, f) = (4, 3.0)$	<b>76.1</b>	$(B, f) = (4, 3.0)$	<b>84.5</b>	$(B, f) = (6, 0.7)$	68.5	$(B, f) = (1, 2.0)$	<b>83.9</b>
SAP	$(B, f) = (6, 0.7)$	70.3	$(B, f) = (6, 2.0)$	83.7	$(B, f) = (1, 3.0)$	69.4	$(B, f) = (6, 1.0)$	74.8
Dropout	$(B, \pi_{dnp}) = (6, 0.3)$	70.1	$(B, \pi_{dnp}) = (6, 0.1)$	83.3	$(B, \pi_{dnp}) = (5, 0.1)$	63.6	$(B, \pi_{dnp}) = (6, 0.1)$	73.7

TABLE 9: AUC-ROC of different attack-detection sampling schemes on cats-and-dogs dataset with against FGSM and MIM attacks. Higher values indicate better detection.

Specifically, let us now model the attacker by seeking perturbation  $\delta$  such that not only classification error, namely cross-entropy, is maximized, but also uncertainty, or variance, of the network output is simultaneously minimized.

Thus, the adversarial objective is updated as

$$\min_{\|\delta\|_\infty < \epsilon} \sum_{k=1}^K y_{\nu,k}^{\text{target}} \log(y_{\nu,k}(\mathbf{x} + \delta)) + \frac{\mu}{R} \sum_{r=1}^R \|\mathbf{y}_\nu^{(r)}(\mathbf{x}_\nu + \delta) - \bar{\mathbf{y}}_\nu(\mathbf{x}_\nu + \delta)\|_2^2. \quad (28)$$

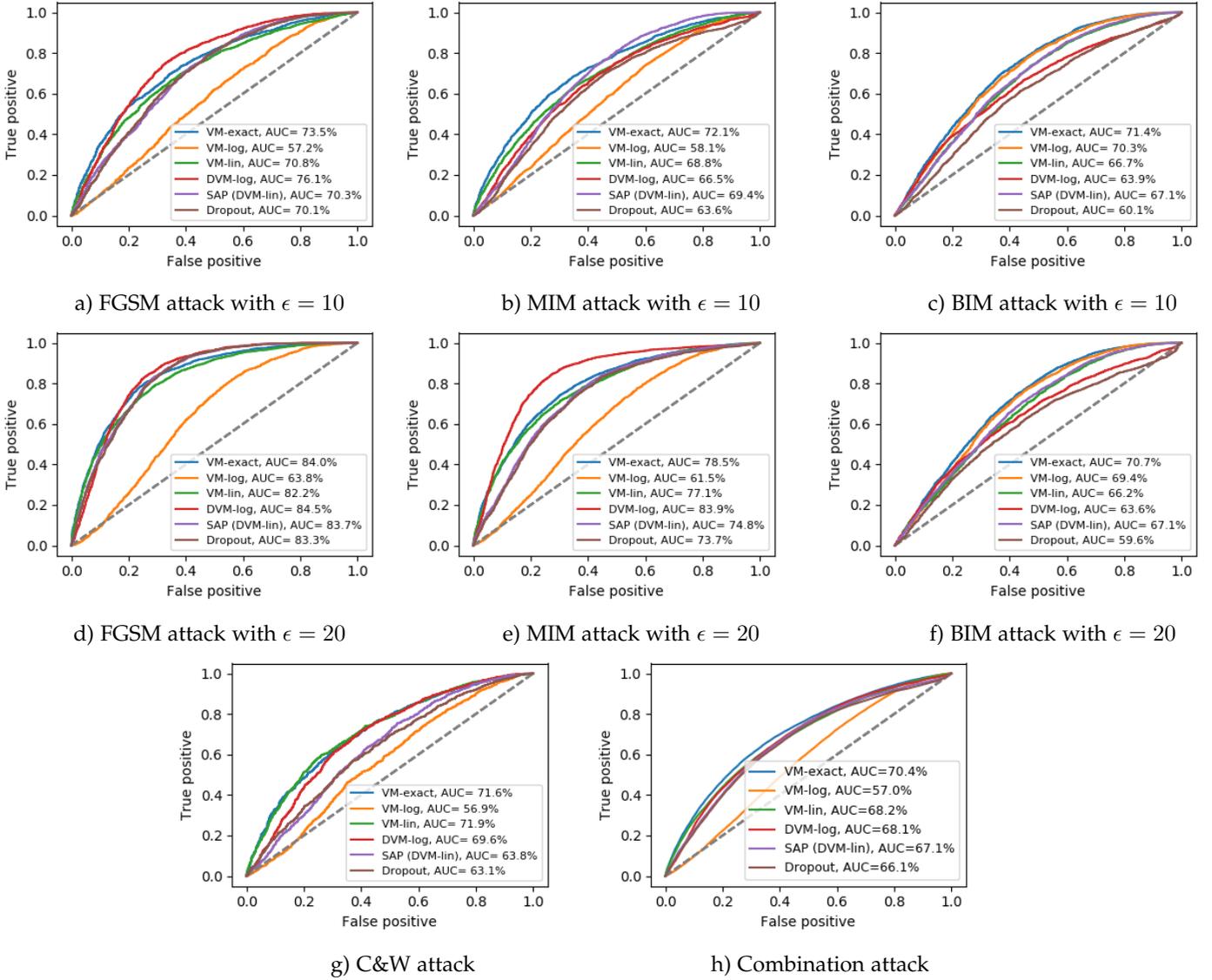


Fig. 4. ROC-curve of different attack-detection sampling schemes on cats-and-dogs dataset against different attacks.

The first term in the objective is the negative cross-entropy between the ground-truth one-hot label of input  $\mathbf{x}_\nu$ , denoted by  $\mathbf{y}_\nu^{\text{target}} = [y_{\nu,1}^{\text{target}}, \dots, y_{\nu,K}^{\text{target}}]$ , and the soft-max output of the deterministic deep neural network  $y_{\nu,k}(\mathbf{x} + \delta)$ , that is with no random sampling unit. The second term is the estimated variance of the detection network output over  $r = 1, \dots, R$  realizations, denoted by  $\mathbf{y}_\nu^{(r)}$ , whose expected value is estimated by its sample average as

$$\bar{\mathbf{y}}_\nu(\mathbf{x}_\nu + \delta) := \frac{1}{R} \sum_{r=1}^R \mathbf{y}_\nu^{(r)}(\mathbf{x}_\nu + \delta) \quad (29)$$

Finally, the scalar  $\mu$  balances the trade-off between maximizing the classification error (chance of a successful attack), and minimizing the output uncertainty, thus mitigating detection. This can potentially lead to lower attack success rate as misclassification is not the sole objective anymore.

Solving (28) analytically is challenging as the sampling probabilities in the random network are in fact a function of the input  $\mathbf{x} + \delta$  as well, making the overall minimization

highly non-convex. In addition, there are several variations of sampling schemes based on different approximations of the output variance. Thus, it is difficult to analytically derive the sampling probabilities for various schemes and substitute them in (28). In this work, we test the performance of the proposed minimum-uncertainty based detection scheme against FGSM attacks targeting the objective in (28), where the attack perturbation  $\delta$  is crafted as the stochastic gradient averaged over  $R = 50$  realization, defined as

$$\delta = \epsilon \cdot \text{sign}(\hat{\mathbf{g}}) \quad (30)$$

with

$$\hat{\mathbf{g}} := \nabla_{\mathbf{x}} \left( \sum_{k=1}^K y_{\nu,k}^{\text{target}} \log(y_{\nu,k}(\mathbf{x})) + \frac{\mu}{R} \sum_{r=1}^R \|\mathbf{y}_\nu^{(r)}(\mathbf{x}) - \bar{\mathbf{y}}_\nu(\mathbf{x})\|_2^2 \right). \quad (31)$$

Similar to Section 6.2, attacks are performed on the cats-and-dogs dataset, on clean images that were correctly classified by the deterministic network. Attack parameters are set as

$\epsilon = 10$  and  $\epsilon = 20$  for values of  $\mu = 0.01, 0.1, 1$  with  $R = 50$ , and the ROC curves are depicted in Fig. 5.

Furthermore, attack success rate corresponding to the percentage of clean images that were correctly classified by the full network and misclassified if perturbed by (31), as well as the AUC-ROC for  $\epsilon = 10$  and varying values of  $\mu$  is reported in Fig. 6. It is interesting to note that in all detection mechanisms, incorporating the variance minimization lowers the attack success rate. However, in DVM-log sampling schemes, where the sampling probabilities are dynamically updated per layer effective in high  $C$  regimes, provides the highest robustness in terms of attack success rate reduction as well as the AUC-ROC. Attack success rate also drops significantly in VM and VM-lin while the AUC-ROCs exhibit slight decrease. In contrast, attack success reduction occurs at higher  $\mu$  values in Dropout and SAP detection schemes, accompanied with a slight increase in the AUC-ROC.

## 7 CONCLUSIONS AND RESEARCH OUTLOOK

Safe and reliable utilization of state-of-the-art CNNs is contingent upon robustifying their performance against adversarial perturbations. To this end, and inspired by Bayesian neural networks, we have investigated attack detection, where through imposing a certain level of randomness we designed the variational distribution to minimize network uncertainty. The premise is that the inherent distance of the adversarial perturbation from the natural-image manifold will cause the overall network uncertainty to exceed that of the clean image, and thus facilitate successful detection. Network uncertainty was expressed as a summation of its layer-wise components, whose exact as well as approximate minimizers have been developed. Links with recent sampling-based approaches have been delineated, along with efficient implementations of the proposed approach. Finally, numerical tests on the CIFAR10 and the cats-and-dogs datasets on deep state-of-the-art CNNs demonstrated the importance of placement as well as tuning of the sampling parameters, which readily translate to improved attack detection.

Among future directions, one can incorporate the novel approach in attack correction schemes based on randomization. This is of particular importance for careful sampling at the initial layers of the network. One can further exploit an ensemble of detection networks, in which sampling units are incorporated at *random* depths. In addition to increased defense strength, the latter introduce a second source of randomness in the defense mechanism, and thus prevent identification by the attackers. Provable detection with performance guarantees is also among future directions.

## 8 APPENDIX

### 8.1 Proof of Proposition 1

Define  $\mathbf{u}_\nu := \mathbf{W}_2\sigma(\mathbf{W}_1\mathbf{x}_\nu)$  and approximate it using the first-order Taylor expansion around  $\bar{\mathbf{u}}_\nu := \mathbb{E}_{q_\theta(\omega)}[\mathbf{u}_\nu]$ , to arrive at

$$\mathbf{y}_\nu \simeq \sigma_{\text{softmax}}(\bar{\mathbf{u}}_\nu) + \nabla\sigma_{\text{softmax}}(\mathbf{u})\Big|_{\mathbf{u}=\bar{\mathbf{u}}_\nu} (\mathbf{u}_\nu - \bar{\mathbf{u}}_\nu) \quad (32)$$

which after taking expectation yields

$$\mathbb{E}_{q_\theta(\omega)}[\mathbf{y}_\nu] \simeq \sigma_{\text{softmax}}(\bar{\mathbf{u}}_\nu). \quad (33)$$

Upon defining the matrix  $\mathbf{H}_1 := \nabla\sigma_{\text{softmax}}(\mathbf{u})\Big|_{\mathbf{u}=\bar{\mathbf{u}}_\nu}$ , and using (32) and (33), we find  $\mathbf{y}_\nu - \mathbb{E}_{q_\theta(\omega)}[\mathbf{y}_\nu] \simeq \mathbf{H}_1(\mathbf{u}_\nu - \bar{\mathbf{u}}_\nu)$  that leads to approximating the variance score as

$$\begin{aligned} & \text{Cov}_{q_\theta(\omega)}[\mathbf{y}_\nu] \\ &= \mathbb{E}_{q_\theta(\omega)}\left[(\mathbf{y}_\nu - \mathbb{E}_{q_\theta(\omega)}[\mathbf{y}_\nu])(\mathbf{y}_\nu - \mathbb{E}_{q_\theta(\omega)}[\mathbf{y}_\nu])^\top\right] \\ &\simeq \mathbb{E}_{q_\theta(\omega)}\left[\mathbf{H}_1(\mathbf{u}_\nu - \bar{\mathbf{u}}_\nu)(\mathbf{u}_\nu - \bar{\mathbf{u}}_\nu)^\top \mathbf{H}_1^\top\right]. \end{aligned}$$

The trace of the latter can be upper bounded by

$$\text{Tr}(\text{Cov}_{q_\theta(\omega)}[\mathbf{y}_\nu]) \leq \lambda_1 \text{Tr}(\text{Cov}_{q_\theta(\omega)}[\mathbf{u}_\nu])$$

where  $\text{Tr}(\text{Cov}_{q_\theta(\omega)}[\mathbf{u}_\nu]) := \text{Tr}(\mathbb{E}_{q_\theta(\omega)}[(\mathbf{u}_\nu - \bar{\mathbf{u}}_\nu)(\mathbf{u}_\nu - \bar{\mathbf{u}}_\nu)^\top])$ , and  $\lambda_1 := \text{Tr}(\mathbf{H}_1^\top \mathbf{H}_1)$  is *deterministic*. Thus, the output variance score is upperbounded by that of the previous layer up to a constant  $\lambda_1$ . Repeating this process of approximating  $\mathbf{u}_\nu$  as a function of  $\mathbf{v}_\nu = \mathbf{W}_1\mathbf{x}_\nu$  by the first-order Taylor expansion around  $\bar{\mathbf{v}}_\nu := \mathbb{E}_{\mathbf{W}_1}[\mathbf{v}_\nu]$ , leads with  $\mathbf{H}_2 := \nabla\sigma(\bar{\mathbf{v}}_\nu)$  to  $\mathbf{u}_\nu$  and its mean compensated approximation

$$\begin{aligned} \mathbf{u}_\nu &\simeq \mathbf{W}_2\sigma(\bar{\mathbf{v}}_\nu) + \mathbf{W}_2\mathbf{H}_2(\mathbf{v}_\nu - \bar{\mathbf{v}}_\nu) \\ \bar{\mathbf{u}}_\nu &\simeq \mathbb{E}_{\mathbf{W}_2}[\mathbf{W}_2\sigma(\bar{\mathbf{v}}_\nu)] \\ \mathbf{u}_\nu - \bar{\mathbf{u}}_\nu &\simeq \mathbb{E}_{\mathbf{W}_2}[\mathbf{W}_2\sigma(\bar{\mathbf{v}}_\nu)] - \mathbf{W}_2\sigma(\bar{\mathbf{v}}_\nu) + \mathbf{W}_2\mathbf{H}_2(\mathbf{v}_\nu - \bar{\mathbf{v}}_\nu). \end{aligned} \quad (34)$$

The latter yields the covariance approximation

$$\begin{aligned} \text{Cov}_{q_\theta(\omega)}[\mathbf{u}_\nu] &\simeq \text{Cov}_{\mathbf{W}_2}[\mathbf{W}_2\sigma(\bar{\mathbf{v}}_\nu)] \\ &+ \mathbb{E}_{\mathbf{W}_2}[\mathbf{W}_2\mathbf{H}_2\mathbb{E}_{\mathbf{W}_1}[(\mathbf{v}_\nu - \bar{\mathbf{v}}_\nu)(\mathbf{v}_\nu - \bar{\mathbf{v}}_\nu)^\top]\mathbf{H}_2^\top\mathbf{W}_2^\top] \\ &+ \mathbb{E}_{\mathbf{W}_2}\left[\mathbf{W}_2\mathbf{H}_2\mathbb{E}_{\mathbf{W}_1}(\mathbf{v}_\nu - \bar{\mathbf{v}}_\nu) \right. \\ &\quad \left. \times \left(\mathbb{E}_{\mathbf{W}_2}[\mathbf{W}_2\sigma(\bar{\mathbf{v}}_\nu)] - \mathbf{W}_2\sigma(\bar{\mathbf{v}}_\nu)\right)^\top\right] \\ &= \text{Cov}_{\mathbf{W}_2}[\mathbf{W}_2\sigma(\bar{\mathbf{v}}_\nu)] + \mathbb{E}_{\mathbf{W}_2}[\mathbf{W}_2\mathbf{H}_2 \text{Cov}_{\mathbf{W}_1}[\mathbf{v}_\nu]\mathbf{H}_2^\top\mathbf{W}_2^\top] \end{aligned}$$

where we have used the independence of random matrices  $\mathbf{W}_1$  and  $\mathbf{W}_2$ , and  $\mathbb{E}_{\mathbf{W}_1}(\mathbf{v}_\nu - \bar{\mathbf{v}}_\nu) = \mathbf{0}$ . Taking the trace and using the inequality  $\text{Tr}(\mathbf{A}\mathbf{B}) \leq \text{Tr}(\mathbf{A})\text{Tr}(\mathbf{B})$  for positive semi-definite matrices  $\mathbf{A}, \mathbf{B} \geq \mathbf{0}$  twice, we arrive after defining  $\lambda_2 := \text{Tr}(\mathbf{H}_2^\top \mathbf{H}_2)$ , at

$$\begin{aligned} & \text{Tr}(\text{Cov}_{q_\theta(\omega)}[\mathbf{u}_\nu]) \\ &\simeq \text{Cov}_{\mathbf{W}_2}[\mathbf{W}_2\sigma(\bar{\mathbf{v}}_\nu)] + \mathbb{E}_{\mathbf{W}_2}[\text{Tr}(\mathbf{W}_2\mathbf{H}_2 \text{Cov}_{\mathbf{W}_1}[\mathbf{v}_\nu]\mathbf{H}_2^\top\mathbf{W}_2^\top)] \\ &\leq \text{Tr}(\text{Cov}_{\mathbf{W}_2}[\mathbf{W}_2\sigma(\bar{\mathbf{v}}_\nu)]) \\ &\quad + \lambda_2 \text{Tr}(\mathbb{E}_{\mathbf{W}_2}[\mathbf{W}_2\mathbf{W}_2^\top])\text{Tr}(\text{Cov}_{\mathbf{W}_1}[\mathbf{v}_\nu]). \end{aligned}$$

Leveraging the last inequality, we can majorize the uncertainty minimization in (14) by that in (17). This is a coupled minimization of layer-wise variance scores  $\text{Tr}(\text{Cov}_{\mathbf{W}_1}[\mathbf{v}_\nu])$  and  $\text{Tr}(\text{Cov}_{\mathbf{W}_2}[\mathbf{W}_2\sigma(\bar{\mathbf{v}}_\nu)])$ , that we solve as follows.

Using  $\mathbf{W}_1^{(TR)}$  along with (4) and (6), we have  $\mathbf{W}_1 = \mathbf{W}_1^{(TR)}\mathbf{S}_1\mathbf{D}_1$ , where  $\mathbf{S}_1 := \text{diag}([z_{1,1}, z_{1,2}, \dots, z_{1,h_1}])$  is the sampling matrix with its pseudo-inverse diagonal mean  $\mathbf{D}_1 := \text{diag}^\dagger(\mathbb{E}_{q(\mathbf{z}; \mathbf{p}_1)}[z_{1,1}, z_{1,2}, \dots, z_{1,h_1}])$ . This implies that  $\bar{\mathbf{v}}_\nu := \mathbb{E}_{\mathbf{W}_1}[\mathbf{v}_\nu] = \mathbf{W}_1^{(TR)}\mathbf{x}_\nu$ , which does not depend on the sampling vector  $\mathbf{p}_1$ . As a result, the minimization in (17) can be readily solved by the proposed subproblems. ■

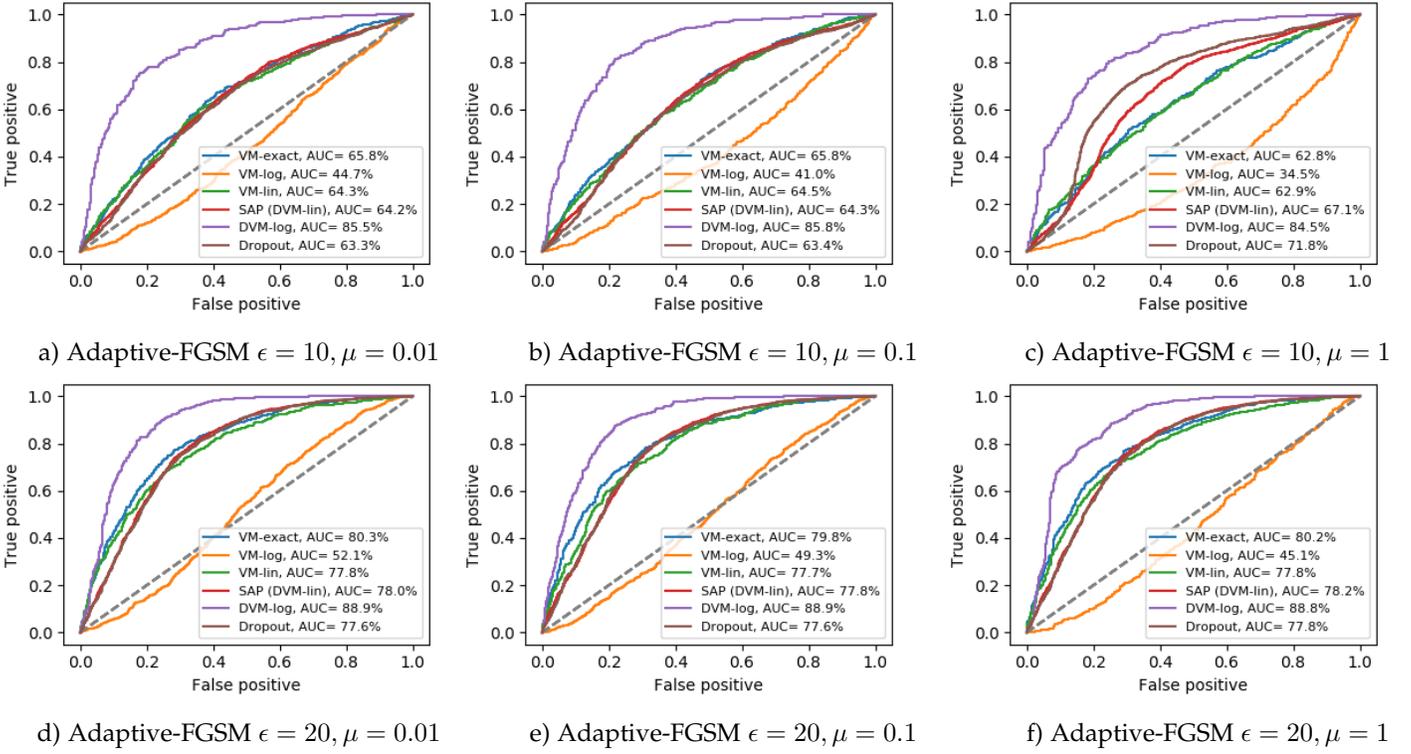


Fig. 5. ROC-curve for different values of  $\mu$  and  $\epsilon$  in adaptive attacks on the cats-and-dogs dataset.

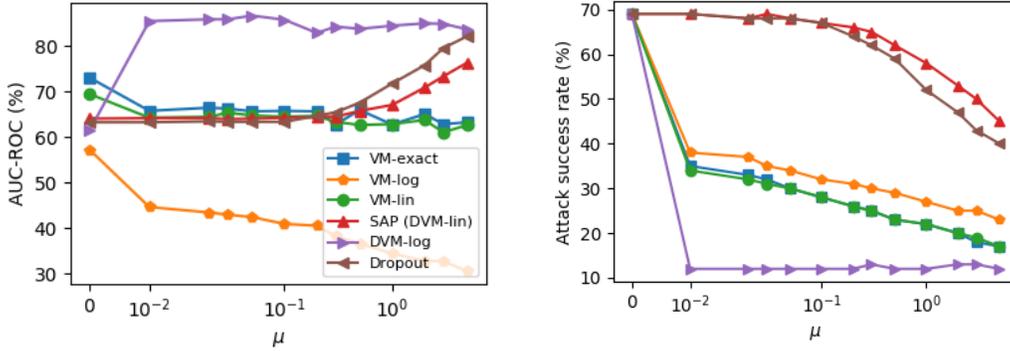


Fig. 6. AUC-ROC and attack success rate against adaptive FGSM attack with  $\epsilon = 10$  for varying values of  $\mu$  for the cats-and-dogs dataset.

Sampling Method	BIM Attack				C&W Attack		Combination Attack	
	$\epsilon = 10$		$\epsilon = 20$		parameters	AUC	parameters	AUC
VM		<b>71.4</b>		<b>70.7</b>				<b>70.4</b>
VM-log	$(B, f) = (1, 3.0)$	70.3	$(B, f) = (1, 3.0)$	69.4	$(B, f) = (2, 0.8)$	56.9	$(B, f) = (2, 0.7)$	57.0
VM-linear		66.7		66.2		<b>71.9</b>		68.2
DVM-log	$(B, f) = (4, 3.0)$	63.9	$(B, f) = (4, 3.0)$	63.6	$(B, f) = (4, 3.0)$	69.6	$(B, f) = (4, 3.0)$	68.1
SAP	$(B, f) = (1, 3.0)$	67.1	$(B, f) = (1, 3.0)$	67.1	$(B, f) = (5, 0.9)$	63.8	$(B, f) = (6, 1.0)$	67.1
Dropout	$(B, \pi_{\text{dnp}}) = (5, 0.6)$	60.1	$(B, \pi_{\text{dnp}}) = (6, 0.1)$	59.6	$(B, \pi_{\text{dnp}}) = (6, 0.5)$	63.1	$(B, \pi_{\text{dnp}}) = (6, 0.1)$	66.1

TABLE 10: AUC-ROC of different attack-detection sampling schemes on cats-and-dogs test set against BIM, C&W, and combination attacks. Higher values indicate better detection.

### 8.2 Proof of Proposition 2

To solve (22), consider the Lagrangian

$$L = \sum_{i=1}^h \frac{\alpha_i}{1 - e^{-Cp_i}} + \rho(\mathbf{1}^\top \mathbf{p} - 1) \quad 0 \leq p_i \leq 1$$

and upon setting its gradient to zero

$$\frac{\partial L}{\partial p_i} = \frac{-C\alpha_i e^{-Cp_i}}{(1 - e^{-Cp_i})^2} + \rho = 0 \quad (35)$$

and introducing the change of variable

$$y_i := \exp(-Cp_i) \quad e^{-C} < y_i < 1$$

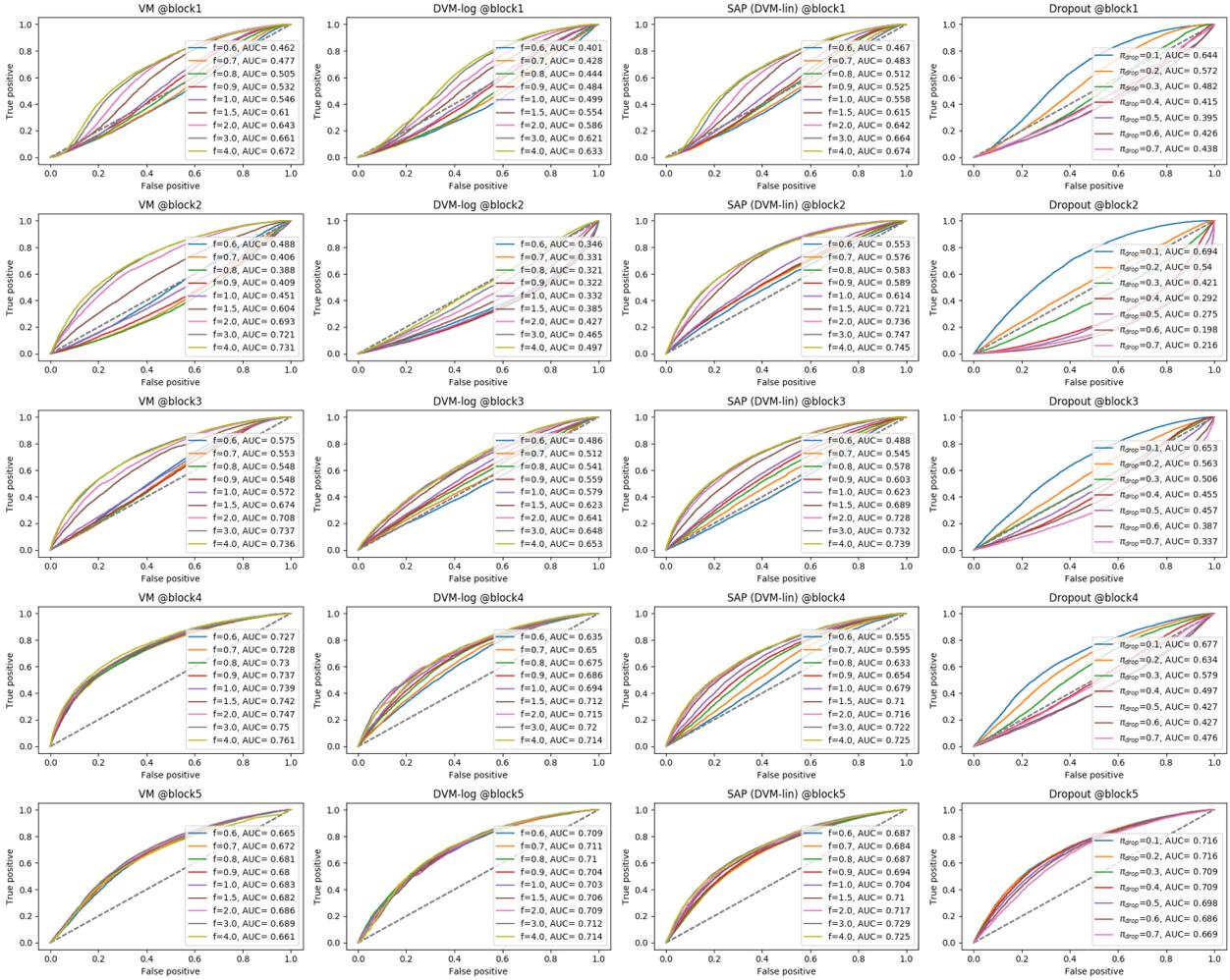


Fig. 7. Performance of different sampling mechanisms at various depths and parameters against combination attack on the CIFAR10 dataset

we find that (35) reduces to  $\rho' y_i^2 - (2\rho' + \alpha_i)y_i + \rho' = 0$ . The feasible root of this quadratic polynomial is

$$y_i = \frac{2\rho' + \alpha_i - \sqrt{(2\rho' + \alpha_i)^2 - 4\rho'^2}}{2\rho'}$$

Using the simplex constraint at the optimal point, we find

$$-\frac{1}{C} \sum_i \ln y_i = 1$$

which after reverting the change of variable, reduces the optimization in (22) to the following root-finding task

$$\sum_i \ln(2\rho' + \alpha_i - \sqrt{(2\rho' + \alpha_i)^2 - 4\rho'^2}) - n \ln(2\rho') + C = 0.$$

This scalar root-finding problem can be solved using bisection that enjoys super-linear convergence rate. ■

### 8.3 Selection of defense parameters

In order to further provide insight on the performance against various selection of  $B$ ,  $f$ , and  $\pi$  parameters, Figs. 7 and 8 illustrate the AUC-ROC for VM-exact, DVM-log, SAP, and uniform dropout against the combination attack. As the plots suggest, uniform dropout reaches its best performance when placed at the last block, whereas higher performance can be obtained by placing carefully-tuned sampling at

units in hidden layers before the last. Furthermore, at a given block  $B$ , VM-exact demonstrates higher robustness for different values of  $f$ ; that is, smaller fluctuation in AUC is observed, whereas other methods are usually more prone to under-performance given sub-optimal parameters.

### REFERENCES

- [1] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *Intl. Conf. on Learning Representations*, San Diego, CA, USA, May 2015.
- [2] J. Redmon and A. Farhadi, “Yolo9000: Better, faster, stronger,” in *IEEE Conf. on Computer Vision and Pattern Recognition*, Honolulu, HI, USA, June 2017, pp. 6517–6525.
- [3] M. Sharif, S. Bhagavatula, L. Bauer, and M. K. Reiter, “Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition,” in *Proc. of Conf. on Computer and Communications Security*, Vienna, Austria, Oct. 2016, pp. 1528–1540.
- [4] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,” in *Advances in Neural Information Processing Systems*, Montréal, Canada, Dec. 2014, pp. 3104–3112.
- [5] J. Engel, C. Resnick, A. Roberts, S. Dieleman, M. Norouzi, D. Eck, and K. Simonyan, “Neural audio synthesis of musical notes with wavenet autoencoders,” in *Intl. Conf. on Machine Learning*, Sydney, Australia, May 2017, pp. 1068–1077.
- [6] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, “Intriguing properties of neural networks,” *Intl. Conf. of Learning Representations*, Banff, Canada, May 2014.
- [7] S. M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, “Universal adversarial perturbations,” in *IEEE Conf. on Computer Vision and Pattern Recognition*, Honolulu, HI, USA, June 2017, pp. 1765–1773.

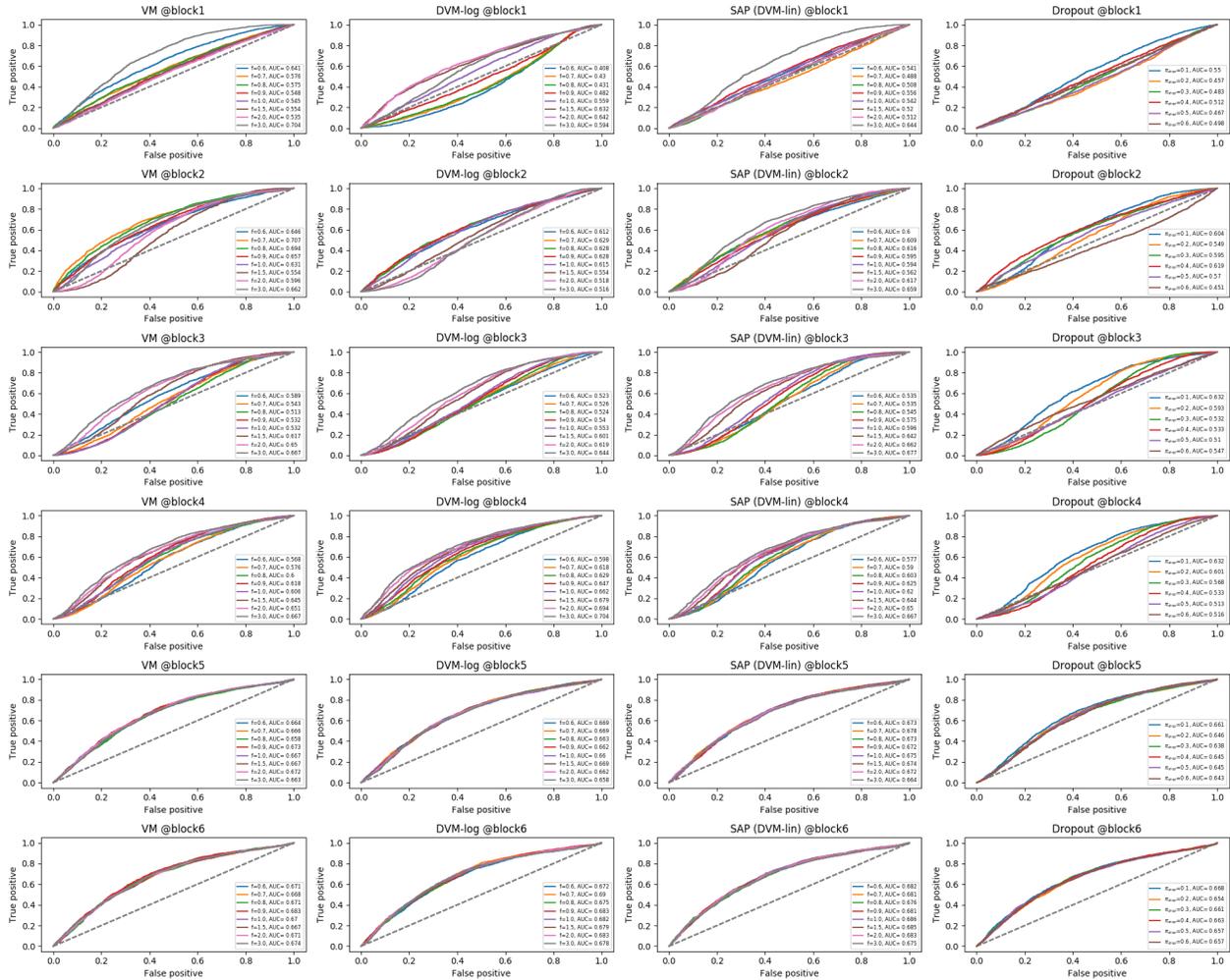


Fig. 8. Performance of different sampling mechanisms at various depths and parameters against combination attack on the cats-and-dogs dataset

[8] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry, "Robustness may be at odds with accuracy," *Intl. Conf. on Learning Representation*, New Orleans, LA, USA, May 2019.

[9] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song, "Robust physical-world attacks on deep learning visual classification," in *IEEE Conf. on Computer Vision and Pattern Recognition*, Salt Lake City, Utah, USA, June 2018, pp. 1625–1634.

[10] G. Zhang, C. Yan, X. Ji, T. Zhang, T. Zhang, and W. Xu, "Dolphinattack: Inaudible voice commands," in *ACM Conf. on Computer and Communications Security*, Dallas, TX, USA, Oct. 2017, pp. 103–117.

[11] G. Goswami, N. Ratha, A. Agarwal, R. Singh, and M. Vatsa, "Unravelling robustness of deep learning based face recognition against adversarial attacks," *AAAI Conf. on Artificial Intelligence*, New Orleans, LA, USA, Feb. 2018.

[12] A. J. Bose and P. Aarabi, "Adversarial attacks on face detectors using neural net based constrained optimization," *IEEE Intl. Workshop on Multimedia Signal Processing*, Vancouver, Canada, Aug. 2018.

[13] N. Carlini and D. Wagner, "Adversarial examples are not easily detected: Bypassing ten detection methods," in *ACM Workshop on Artificial Intelligence and Security*, Dallas, TX, USA, Oct. 2017, pp. 3–14.

[14] —, "Towards evaluating the robustness of neural networks," in *IEEE Symposium on Security and Privacy*, San Jose, CA, USA, May 2017, pp. 39–57.

[15] N. Papernot, P. McDaniel, and I. Goodfellow, "Transferability in machine learning: From phenomena to black-box attacks using adversarial samples," *Technical Report, arXiv preprint arXiv:1605.07277*, 2016.

[16] H. Zhang, H. Chen, Z. Song, D. Boning, I. S. Dhillon, and C.-J. Hsieh, "The limitations of adversarial training and the blind-spot attack," *Intl. Conf. on Learning Representations*, New Orleans, LA, USA, May 2019.

[17] S. Liu, P.-Y. Chen, X. Chen, and M. Hong, "Signsgd via zeroth-order oracle," *Intl. Conf. on Learning Representations*, New Orleans, LA, USA, May 2019.

[18] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, "The limitations of deep learning in adversarial settings," in *IEEE European Symposium on Security and Privacy*, Saarbrücken, Germany, March 2016, pp. 372–387.

[19] S. Gu and L. Rigazio, "Towards deep neural network architectures robust to adversarial examples," *Intl. Conf. on Learning Representations*, San Diego, CA, USA, May 2015.

[20] J. H. Metzen, T. Genewein, V. Fischer, and B. Bischoff, "On detecting adversarial perturbations," *Intl. Conf. on Learning Representations*, Toulon, France, April 2017.

[21] J. Lu, T. Issaranon, and D. Forsyth, "SafetyNet: Detecting and rejecting adversarial examples robustly," in *IEEE Intl. Conf. on Computer Vision*, Venice, Italy, Oct. 2017, pp. 446–454.

[22] R. Feinman, R. R. Curtin, S. Shintre, and A. B. Gardner, "Detecting adversarial samples from artifacts," *arXiv preprint arXiv:1703.00410*, 2017.

[23] J. M. Cohen, E. Rosenfeld, and J. Z. Kolter, "Certified adversarial robustness via randomized smoothing," in *Intl. Conf. on Machine Learning (ICML)*, Long Beach, CA, June 2019.

[24] C. Guo, M. Rana, M. Cisse, and L. Van Der Maaten, "Countering adversarial images using input transformations," *Intl. Conf. on Learning Representations*, Vancouver, Canada, May 2018.

[25] S. Gopalakrishnan, Z. Marzi, U. Madhow, and R. Pedarsani, "Combating adversarial attacks using sparse representations," *arXiv preprint arXiv:1803.03880*, 2018.

[26] T. Miyato, S. Maeda, S. Ishii, and M. Koyama, "Virtual adversarial training: A regularization method for supervised and semi-

- supervised learning," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2018.
- [27] L. Schmidt, S. Santurkar, D. Tsipras, K. Talwar, and A. Madry, "Adversarially robust generalization requires more data," *Advances in Neural Information Processing Systems*, Montreal, Canada, Dec. 2018.
- [28] A. Sinha, H. Namkoong, and J. Duchi, "Certifying some distributional robustness with principled adversarial training," *Intl. Conf. of Learning Representations*, Vancouver, Canada, April 2018.
- [29] E. Wong and J. Z. Kolter, "Provable defenses against adversarial examples via the convex outer adversarial polytope," in *Intl. Conf. on Machine Learning (ICML)*, Stockholm, Sweden, June 2018.
- [30] H. Zhang, H. Chen, C. Xiao, B. Li, D. Boning, and C.-J. Hsieh, "Towards stable and efficient training of verifiably robust neural networks," *arXiv preprint arXiv:1906.06316*, 2019.
- [31] Y. Guo, C. Zhang, C. Zhang, and Y. Chen, "Sparse DNNs with improved adversarial robustness," in *Advances in Neural Information Processing Systems*, Montreal, Canada, Dec. 2018, pp. 240–249.
- [32] A. N. Gomez, I. Zhang, K. Swersky, Y. Gal, and G. E. Hinton, "Targetted dropout," in *Advances in Neural Information Processing Workshop Track*, Montreal, Canada, USA, Dec. 2018.
- [33] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," *Intl. Conf. on Learning Representations*, Toulon, France, April 2017.
- [34] H. Gouk, E. Frank, B. Pfahringer, and M. Cree, "Regularisation of neural networks by enforcing Lipschitz continuity," *arXiv preprint arXiv:1804.04368*, 2018.
- [35] D. Hendrycks and K. Gimpel, "A baseline for detecting misclassified and out-of-distribution examples in neural networks," *Proc. of Intl. Conf. on Learning Representations*, Toulon, France, May 2017.
- [36] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, "On calibration of modern neural networks," pp. 1321–1330, Sydney, Australia, Aug. 2017.
- [37] A. Malinin and M. Gales, "Predictive uncertainty estimation via prior networks," in *Advances in Neural Information Processing Systems*, Montreal, Canada, Dec. 2018, pp. 7047–7058.
- [38] L. Smith and Y. Gal, "Understanding measures of uncertainty for adversarial example detection," *Conf. on Uncertainty in Artificial Intelligence*, Monterey, CA, USA, Aug. 2018.
- [39] A. Kendall and Y. Gal, "What uncertainties do we need in bayesian deep learning for computer vision?" in *Advances in Neural Information Processing Systems*, Long Beach, CA, Dec. 2017, pp. 5574–5584.
- [40] F. Sheikholeslami, S. Jain, and G. B. Giannakis, "Efficient randomized defense against adversarial attacks in deep convolutional networks," in *Proc. of Intl. Conf. on Acoustics, Speech, and Signal Processing*, Brighton, UK, May 12–17, 2019.
- [41] Y. Gal and Z. Ghahramani, "Bayesian convolutional neural networks with Bernoulli approximate variational inference," in *Intl. Conf. on Learning Representations Workshop Track*, San Juan, Puerto Rico, May 2016.
- [42] —, "A theoretically grounded application of dropout in recurrent neural networks," in *Advances in Neural Information Processing Systems*, Barcelona, Spain, Dec. 2016, pp. 1019–1027.
- [43] L. Cardelli, M. Kwiatkowska, L. Laurenti, N. Paoletti, A. Patane, and M. Wicker, "Statistical guarantees for the robustness of Bayesian neural networks," *arXiv preprint arXiv:1903.01980*, 2019.
- [44] G. S. Dhillon, K. Azzadenesheli, Z. C. Lipton, J. Bernstein, J. Kossaifi, A. Khanna, and A. Anandkumar, "Stochastic activation pruning for robust adversarial defense," *Intl. Conf. on Learning Representationa*, Vancouver, Canada, April 2018.
- [45] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [46] C. Bishop, *Pattern Recognition and Machine Learning*, Springer, New York, USA, 2006.
- [47] R. M. Neal, *Bayesian Learning for Neural Networks*. Springer Science & Business Media, New York, USA, 2012.
- [48] D. Garcia-Gasulla, F. Pares, A. Vilalta, J. Moreno, E. Ayguade, J. Labarta, U. Cortes, and T. Suzumura, "On the behavior of convolutional nets for feature extraction," *Journal of Artificial Intelligence Research*, vol. 61, pp. 563–592, 2018.
- [49] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conf. on Computer Vision and Pattern Recognition*, Las Vegas Valley, NV, USA, June 2016, pp. 770–778.
- [50] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *Intl. Conf. on Learning Representations*, San Diego, CA, USA, May 2015.

- [51] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," *Intl. Conf. on Learning Representations*, Toulon, France, April 2017.
- [52] Y. Dong, F. Liao, T. Pang, H. Su, J. Zhu, X. Hu, and J. Li, "Boosting adversarial attacks with momentum," in *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, Salt Lake City, UT, USA, June 2018, pp. 9185–9193.



**Fatemeh Sheikholeslami** (S'13) received the B.Sc. and M.Sc. degrees in Electrical Engineering from the University of Tehran and Sharif University of Technology, Tehran, Iran, in 2010 and 2012, respectively. Since September 2013, she has been working toward the Ph.D. degree with the Department of Electrical and Computer Engineering, University of Minnesota, MN, USA. Her research interests include Machine Learning and Network Science.



**Swayambhoo Jain** (S'13) received the Bachelor of Technology degree in Electronics and Communication Engineering from the National Institute of Technology, Calicut, India in 2007. He was a Research and Development Engineer with IBM Systems and Technology Group, Bangalore, India, from 2007 to 2010. He received the M.Sc. degree and Ph.D. degree in Electrical Engineering from the University of Minnesota, MN, USA, in 2012 and 2017, respectively. Since 2017 he has been working as a Researcher at Technicolor Artificial Intelligence Lab in Palo Alto, CA, USA. His general research interests lie in Machine Learning, Signal Processing, and High Dimensional Statistical Inference. His current research focus is on issues related to efficiency and robustness in Deep Learning.



**G. B. Giannakis** (Fellow'97) received his Diploma in Electrical Engr. from the Ntl. Tech. Univ. of Athens, Greece, 1981. From 1982 to 1986 he was with the Univ. of Southern California (USC), where he received his MSc. in Electrical Engineering, 1983, MSc. in Mathematics, 1986, and Ph.D. in Electrical Engr., 1986. He was with the University of Virginia from 1987 to 1998, and since 1999 he has been a professor with the Univ. of Minnesota, where he holds an Endowed Chair in Wireless Telecommunications, a University of Minnesota McKnight Presidential Chair in ECE, and serves as director of the Digital Technology Center.

His general interests span the areas of communications, networking and statistical learning - subjects on which he has published more than 450 journal papers, 750 conference papers, 25 book chapters, two edited books and two research monographs (h-index 140). Current research focuses on Data Science, Internet of Things, and Network Science with applications to social, brain, and power networks with renewables. He is the (co-) inventor of 32 patents issued, and the (co-) recipient of 9 best journal paper awards from the IEEE Signal Processing (SP) and Communications Societies, including the G. Marconi Prize Paper Award in Wireless Communications. He also received Technical Achievement Awards from the SP Society (2000), from EURASIP (2005), a Young Faculty Teaching Award, the G. W. Taylor Award for Distinguished Research from the University of Minnesota, and the IEEE Fourier Technical Field Award (inaugural recipient in 2015). He is a Fellow of EURASIP, and has served the IEEE in a number of posts, including that of a Distinguished Lecturer for the IEEE-SPS.