

# No Surprises: Training Robust Lung Nodule Detection for Low-Dose CT Scans by Augmenting with Adversarial Attacks

Siqi Liu, Arnaud Arindra Adiyoso Setio, Florin C. Ghesu, Eli Gibson, Sasa Grbic, Bogdan Georgescu, Dorin Comaniciu \*

October 30, 2020

## Abstract

Detecting malignant pulmonary nodules at an early stage can allow medical interventions which may increase the survival rate of lung cancer patients. Using computer vision techniques to detect nodules can improve the sensitivity and the speed of interpreting chest CT for lung cancer screening. Many studies have used CNNs to detect nodule candidates. Though such approaches have been shown to outperform the conventional image processing based methods regarding the detection accuracy, CNNs are also known to be limited to generalize on under-represented samples in the training set and prone to imperceptible noise perturbations. Such limitations can not be easily addressed by scaling up the dataset or the models. In this work, we propose to add adversarial synthetic nodules and adversarial attack samples to the training data to improve the generalization and the robustness of the lung nodule detection systems. To generate hard examples of nodules from a differentiable nodule synthesizer, we use projected gradient descent (PGD) to search the latent code within a bounded neighbourhood that would generate nodules to decrease the detector response. To make the network more robust to unanticipated noise perturbations, we use PGD to search for noise patterns that can trigger the network to give over-confident mistakes. By evaluating on two different benchmark datasets containing consensus annotations from three radiologists, we show that the proposed techniques can improve the detection performance on real CT data. To understand the limitations of both the conventional

---

\*Siqi Liu, Florin C. Ghesu, Eli Gibson, Sasa Grbic, Bogdan Georgescu and Dorin Comaniciu are with Digital Technology & Innovation, Siemens Healthineers, Princeton, NJ, USA. (siqi.liu@siemens-healthineers.com)

Arnaud Arindra Adiyoso Setio is with Digital Technology & Innovation, Siemens Healthineers, Erlangen, Germany.

The authors thank the National Cancer Institute for access to NCI's data collected by the National Lung Screening Trial (NLST). The statements contained herein are solely those of the authors and do not represent or imply concurrence or endorsement by NCI.

This manuscript has been published on IEEE Trans. on Medical Imaging.

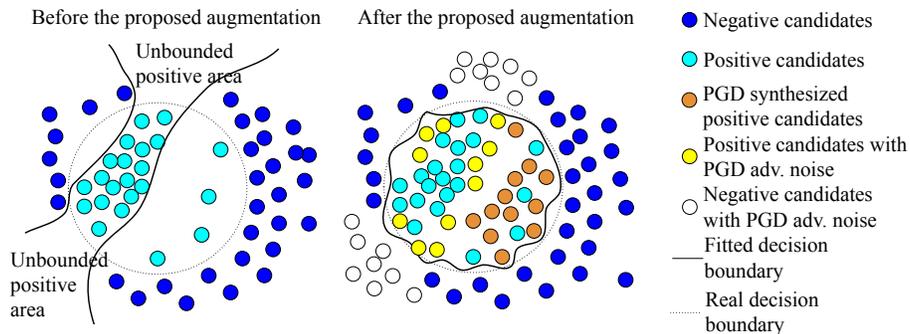


Figure 1: A conceptual illustration of the proposed training scheme. Pulmonary nodules in chest CTs follow a long-tail distribution with rare and hard nodules under-represented. ReLU networks tend to form open decision boundaries which leave the risk for the network to be activated by arbitrary noise [5]. We propose adversarial augmentation methods to efficiently search for both hard synthetic nodules and adversarial samples that can improve the robustness of the network.

networks and the proposed augmented networks, we also perform stress-tests on the false positive reduction networks by feeding different types of artificially produced patches. We show that the augmented networks are more robust to both under-represented nodules as well as resistant to noise perturbations.

## 1 Introduction

Lung cancer is the leading cause of all cancer deaths [1]. Detecting malignant pulmonary nodules at an early stage can allow medical interventions which may increase the survival rate of lung cancer patients. Early-stage cancer generally manifests in the form of pulmonary nodules which are defined as rounded opacity, well or poorly defined, measuring up to 30mm in diameter [2]. Based on the findings of the National Lung Screening Trial (NLST), the U.S. Centers for Medicare and Medicaid Services (CMS) approved screening for lung cancer of high-risk subjects to be fully reimbursed by insurance companies. In NLST the low-dose screening test involved an approximate dose of 2 mSv, whereas full-chest CT scanning that was the major diagnostic study used to follow up nodules, involved a dose of about 8 mSv [3]. The NELSON trial also reported reduced 10 year lung-cancer mortality with CT screening with a randomized trial involving 15789 patients [4]. However, given the sizeable eligible screening population (8.6 million in the US) and the time cost of interpreting 3D chest CT, it substantially increases the efforts for radiologists.

Motivated by the LUNA16 challenge [6], many studies have attempted to automate the detection of pulmonary nodules using machine learning, in particular, deep convolutional neural networks (CNN) in order to assist the radiologists in the lung screening workflow [7, 8, 9, 10, 11, 12]. Following the coarse-to-fine strategy, the majority of the deep learning-based nodule detection methods are implemented as a two-stage system:

(1) a candidate generation network with a large field of view is first trained to output initial detection results with high sensitivity at the cost of low specificity; (2) a false positive reduction (FPR) network is then trained to re-evaluate the confidence of each candidate.

Though many show CNNs can improve both the sensitivity and the specificity comparing to the previous image processing based CAD systems, CNNs can suffer from a few challenges, which we argue cannot be addressed by simply adding more training data or hyper-parameter tuning. First, the observer variability among radiologists is known to be high. For example, only 928 out of 2669 suspected findings from the LIDC-IDRI study are agreed as nodules ( $\geq 3\text{mm}$ ) by all the four radiologists [13]. Such variability can be caused by factors such as the vague definition of pulmonary nodules, the imbalanced level of expertise among radiologists or the insufficient information provided by chest CT, etc. Second, the detection networks tend to miss nodules that are under-represented in the training set, such as the small ground-glass nodules, irregularly shaped nodules or nodules appearing in under-represented contexts. Because only 3.6% of the screening population have biopsy-proven malignant nodules [14], such malignant nodules can also be under-represented in the training data. Third, neural networks are known to be prone to unexpected image distortions [15]. Such distortions can happen in the real-world low-dose CT imaging though they are rare in both the training and the benchmark datasets. As we show later in this paper, even simple noise patterns can determine an under-augmented nodule detector to giving positive responses. Under- or over-detecting nodules caused by such unanticipated distortions can pose the potential risk of distracting and biasing the radiologists. Therefore, besides achieving overall high sensitivity and a low number of false positives on clean benchmark datasets, a nodule detection system is also expected to (1) be capable of detecting under-represented nodules that are rare in both the training and benchmark datasets (2) be robust to unanticipated noise and distortions in the real-world images.

Motivated by the reasons above, we propose to augment the training set of lung nodule detection by adversarially attacking a pre-trained false positive reduction network with both hard synthetic nodules as well as noise image perturbations. The concept is illustrated in Fig. 1. First, we use projected gradient descent (PGD) [16] to search for the adversarial samples that can determine a trained false positive reduction network into outputting over-confident wrong predictions. These searched patches are then added to the training patches to augment the detector to be more robust to both under-represented nodules and unanticipated image distortions. PGD is used for searching for three types of adversarial augmentation patches: (1) latent codes to sample hard synthetic nodules that the detector fails to detect; (2) perturbation noise that can make the nodule detector fail to detect; (3) noise patterns that can easily determine the nodule detector to giving false-positive findings; To evaluate the proposed methods, we train a baseline nodule detector following the general 2-stage framework using a large-scale training dataset. The adversarial patches are then generated by attacking the baseline false positive reduction (FPR) network and are used for augmenting the FPR network. By evaluating on two different benchmark datasets, we show the proposed techniques can improve the detection performance on clean benchmark data. Using the same techniques, we also generate adversarial samples to stress-test the trained false positive reduction networks. We show that the augmented networks are

more robust to both hard nodules and noise perturbations.

## 2 Related Work

### 2.1 Deep learning based nodule detection

As one of the most popular applications of computer-aided diagnosis systems, many studies have been dedicated to using image processing and machine learning algorithms to detect lung nodules [17]. The majority of the nodule detection framework generate candidates first either with an image processing pipeline or a fully convolutional neural network. Then a separate classifier is trained to reduce false positives based on the input 3D CT patches centered at the candidate locations. Most of the recent works were developed based on the LUNA challenge [6] which acquired its data from the LIDC-IDRI dataset [13]. Though the annotation process of the LIDC-IDRI dataset has been well documented and is considered reliable, the quantity and diversity of the LIDC-IDRI dataset are highly limited. Besides the LUNA challenge, there have been no benchmarks reported with known statistics. Though the metrics computed from the FROC curves are suitable for reporting the detection performance on a given benchmark dataset, it is not often thoroughly investigated that how robust such detection systems would perform on the rare cases as well as noise perturbations. Our work shows that the conventional CNNs trained without adversarial augmentation would generally fail to recognize rare nodules as well as prone to image noise. For a more comprehensive review of the deep learning based lung nodule detection systems, we would refer our readers to [18, 17].

### 2.2 Data synthesis based augmentation in medical image analysis

Inspired by the recent advances in generative models, there have been increasing interests in synthesizing objects in medical images to augment the existing training set for better diversity [19]. Many recent studies proposed to use generative networks to synthesize lung nodules to improve the performance of diverse lung nodule related applications [20, 21, 22, 23, 24, 25, 26, 27]. Most learning based nodule synthesis methods start with training a generative network to map low dimensional latent codes to realistic lung nodules in chest CT using either variational auto-encoder (VAE) or Generative adversarial networks (GAN). Latent codes are sampled from a predefined prior distribution randomly to synthesize nodules resembling the real ones. These synthetic nodules are blended into the original image contexts by either formulating the training task as either image inpainting [22] or using an extra context-blending network [21]. In [21], authors use both the discriminator error and the classification error to select only the hard synthetic cases to be added to the augmented dataset. We show that such sampling strategies can be inefficient. The majority of the synthetic samples would add little values since they can be successfully recognized by a network that is trained on a large-scale dataset. However, hard samples can be drawn from a synthesizer without exhaustive search if the latent codes are optimized to increase the training loss of a trained network. In [28], authors showed adversarial sampling can

help network generalize better on multi-class classification problems.

### **2.3 Over-confident neural networks and adversarial training**

To build robust computer-aided diagnosis systems that are robust to out of distribution (OOD) samples, one can train the network to estimate the decision uncertainty and reject the samples when the estimated uncertainty is high [29, 30]. Though we also use the beta distribution in our work for uncertainty estimation [29], we show that the uncertainty estimation techniques alone would be insufficient to make the network robust to avoid over-confident decisions on OOD samples. In [16], it is argued that ReLU activated neural networks would always have open decision boundaries which leave the risk of high responses for unseen OOD samples. In another paper, it is argued that batch normalization is also a cause of the adversarial vulnerability [31]. Such network vulnerability is hard to be reflected by the clean medical image benchmark datasets. In [32, 5], it is proposed to use PGD [16] to search for the adversarial augmentation cases from uniform noise or permuted input patches to augment the clean training dataset. We use similar techniques to adversarially sample both hard positive and hard negative nodule samples to enhance the adversarial robustness of the nodule detection networks. Though it was suggested that the adversarially trained networks can generalize slightly worse on clean data [33, 34, 35, 36], we believe such robustness is still vital for real-world medical AI applications.

### **2.4 Adversarial robustness of medical image analysis systems**

The vulnerability of CNN against adversarial noise also poses potential risks for deploying the computer-aided diagnosis systems in real clinics as investigated by some recent studies [37, 38, 39]. Some early studies also attempted to defend the networks from adversarial noise using different types of data augmentation, such as using geometric transformation [40] or adding Poisson noise [41]. Authors of [42] also propose to use the model ensemble to improve the model robustness of nodule malignancy prediction network. Given the fact that defending against adversarial samples is a challenging task, [43, 44] also proposed to analyze the neural network feature distributions to detect adversarial samples.

## **3 Methods**

### **3.1 Baseline Detection Architectures**

Similar to many new deep learning based nodule detection frameworks, our baseline framework consists of a candidate generation (CG) module and a false positive reduction (FPR) module as shown in Fig. 3. The candidate generation module is trained to achieve high sensitivity via over-detecting nodule candidates. We use three identical 3D ResUNets [45] as the CG backbone networks without weight sharing. The first CG network is first trained to output 3D heatmaps with the nodule centers represented by 3D Gaussian blobs with the same sizes (3D Blob All Nodules). We then fine-tune the

first CG with only the ground glass candidates and part-solid candidates since they are under-represented in the training set (3D Blob Ground Glass Nodules). The candidates are derived with non-maximum suppression (NMS) on the fusion heatmap obtained by taking the element-wise maxima of the two network output heatmaps. We also fine-tune the first CG network by adding a 3D region proposal network (RPN) head [46] to outputting 3D bounding boxes (3D RPN Head). We observed that even though it was hard to improve the sensitivity of the standalone 3D RPN based CG alone (72.00% and 73.07% sensitivity in both reported benchmarks), some of the true positive findings are complementary to the blob based CG. By merging the 3D RPN and blob CG candidates, we improved the blob CG sensitivities from 98.00% and 93.07% to 100% and 97.69% when having 100 candidates per scan. The final candidates of the system are obtained by taking the union of the blob candidates and the 3D RPN bounding box candidates.

The false-positive reduction module is then trained to re-evaluate the candidates and prune the false-positive findings based on the classification confidence. It is built with a DenseUNet network pre-trained with nodule segmentation. We add shallow classifier layers on top of it to derive the FPR confidence scores. The network is trained using  $64^3$  patches with a resolution of 0.625<sup>3</sup>mm. We train all the CG and FPR networks using the Adam optimizer [47] with the initial learning rate 0.001.

We trained the CG framework first and froze it before performing the analysis presented in this work. For the brevity of this paper, we demonstrate the proposed techniques only to improve the FPR while assuming the CG networks are trained and frozen. However, the same techniques can also be used for improving CG networks.

### 3.2 Hard-Sample Synthesis with PGD Sampling

We train a nodule synthesizer  $f_{generator}$  that can be controlled by the latent code sampled from a prior distribution. We implement the  $f_{generator}$  with a 3D convolutional variational encoder. We extract the nodules out of the CT context with the manually annotated nodule segmentation. The boundary of the nodule segmentation is blurred with a distance transform. As shown in Fig. 4, we firstly map the cropped 3D nodules to an encoding space using the encoder network  $f_{encoder}$ , then the variational encoding is reconstructed back to the nodules in chest CT. We jointly train a WGAN-GP discriminator [48] with spectral normalization [49] to enforce the generator to add high-frequency details to mimic the real nodules in CT. The data flow can be summarized as

$$\mu_i, \sigma_i = f_{encoder}(x_i^{nodule}) \quad (1)$$

$$z_i^{nodule} \sim \mathcal{N}(\mu_i, \sigma_i^2) \quad (2)$$

$$\tilde{x}_i^{nodule} = f_{generator}(z_i^{nodule}) \quad (3)$$

$$d_i^{fake}, d_i^{real} = f_{discriminator}(\tilde{x}_i^{nodule}, x_i^{nodule}) \quad (4)$$

Here  $d_i^{fake}$  and  $d_i^{real}$  the discriminator output for the fake and real samples. The training objective of the nodule synthesizer can be summarized as

$$L_{discriminator} = L_{WGAN-GP}(d_i^{fake}, d_i^{real}) \quad (5)$$

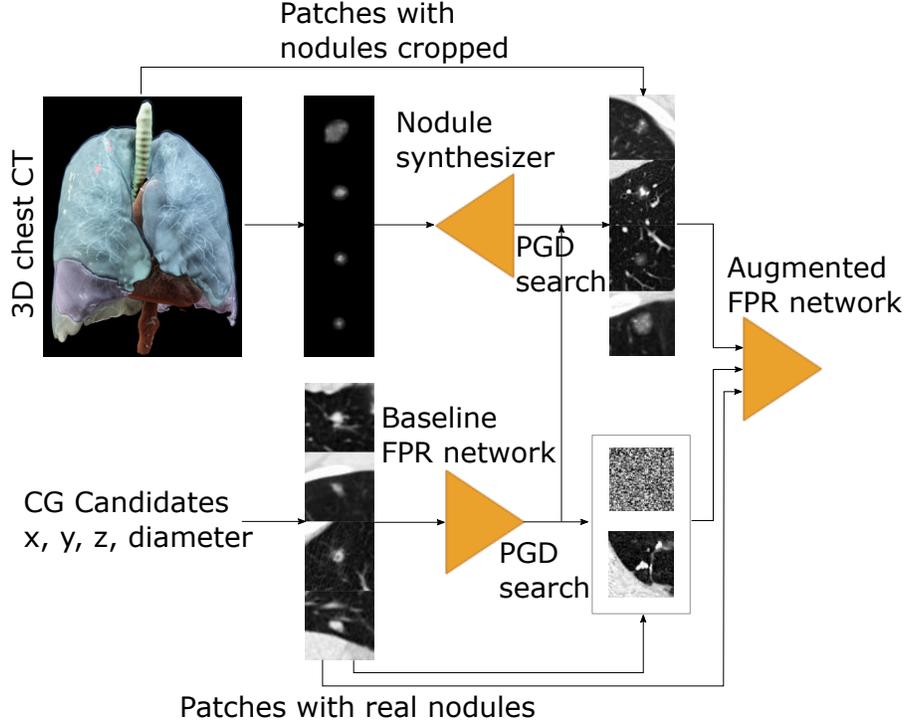


Figure 2: The data-flow illustration of the proposed adversarial augmentation framework for enhancing the false positive reduction (FPR) network in a nodule detection pipeline.

$$\begin{aligned}
 L_{encoder} + L_{generator} = & |\tilde{x}_i^{nodule} - x_i^{nodule}| + \\
 & \lambda_1 D_{KL}(\mathcal{N}(\mu_i, \sigma_i^2) || \mathcal{N}(0, 1)) - \\
 & \lambda_2 L_{WGAN-GP}(d_i^{fake}, d_i^{real})
 \end{aligned} \tag{6}$$

where  $D_{KL}(\mathcal{N}(\mu_i, \sigma_i^2) || \mathcal{N}(0, 1))$  optimizes the probability distribution parameters  $\mu$  and  $\theta$  to closely resemble that of  $\mathcal{N}(0, 1)$ .  $\lambda_2 L_{WGAN-GP}(d_i^{fake}, d_i^{real})$  is the wasserstein GAN discriminator loss regularized by the gradient penalty defined in [48]. In our experiments, we set  $\lambda_1 = 10^{-5}$  and  $\lambda_2 = 0.1$ .

Once the synthesizer is trained, we discard both the encoder network and the discriminator. Only the generator network is kept for sampling synthetic nodules. Random nodules can be sampled by feeding a code to the trained generator  $f_{generator}(z_i^{nodule} \sim \mathcal{N}(\mu_i, \sigma_i^2))$ . The synthesized nodule can be fused to a random background chest CT patch  $x_i$  and then fed to a trained FPR classifier  $f_{FPR}$ . For the rest of the paper, we define the differentiable fusion of synthetic nodule and the background as  $\tilde{x}_i^{nodule} \oplus x_i = \tilde{x}_i^{nodule} * m_i + x_i * (1 - m_i)$  where  $m_i$  is a binary mask obtained by thresholding  $\tilde{x}_i^{nodule}$ . Though it is feasible to add another training stage as described in [21] to further blend the generated nodule into its context, we found it non-critical for

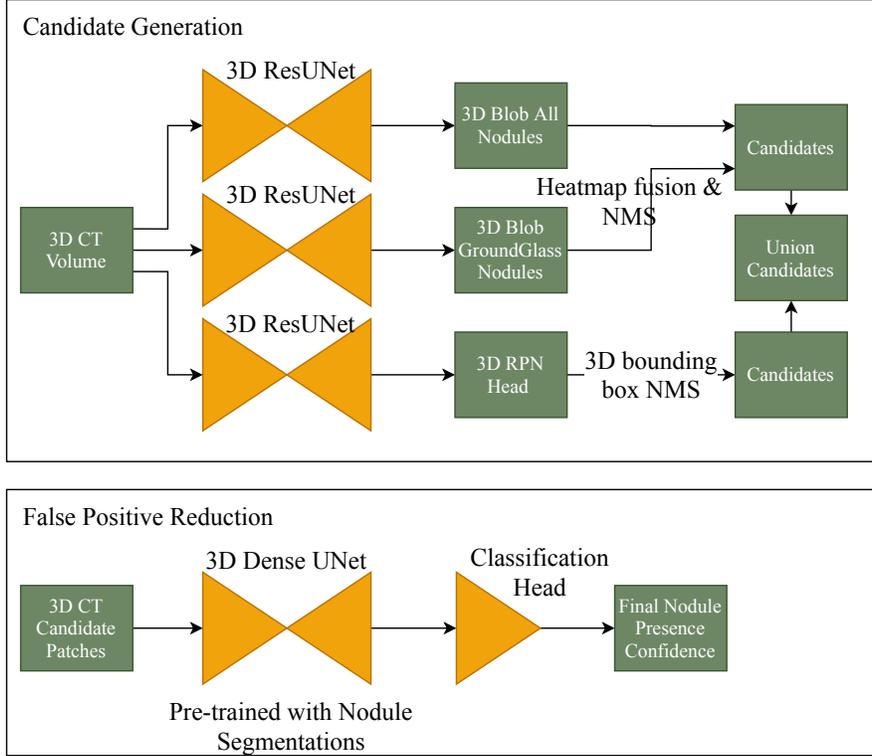


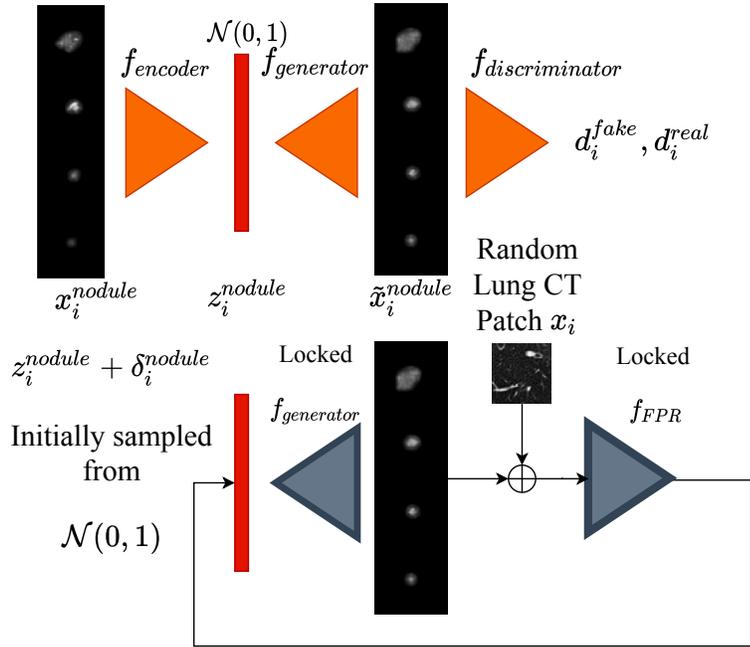
Figure 3: The baseline two stage nodule detection framework used in this work.

the sake of improving the nodule detection in practice. It is inefficient to draw hard-cases directly by randomly sampling from the prior because most of the cases close to the mean have already been learned by the nodule false positive reduction network  $f_{FPR}$ . So instead of randomly sampling the encoding of nodules, we use the projected gradient descent (PGD) as originally used for generating adversarial attacks [16] to sample hard nodules. For each sampling, we initialize the encoding from the standard normal distribution  $z_i^{nodule} \sim \mathcal{N}(0, 1)$  and randomly initialize a perturbation vector  $\delta_i^{nodule}$  to explore the neighbourhood  $\mathcal{S}$  of  $z_i^{nodule}$  within a bounded radius.  $\delta_i^{nodule}$  is updated by PGD to maximize the  $L_{FPR}$  as

$$e_i = f_{FPR}(f_{generator}(z_i^{nodule} + \delta_i^{nodule}) \oplus x_i) \quad (7)$$

$$\arg \max_{\|\delta\| \leq \epsilon} L_{FPR}(e_i, 1) \quad (8)$$

Here,  $\oplus$  is the fusion operator that blends the synthetic nodule into the CT context patch  $x_i$ . As an alternative to the combination of the sigmoid activation and the binary cross-entropy loss, We also use the beta distribution as in [29] to measure the classification uncertainty in our experiments. With the beta distribution output, the FPR



PGD search for the code that could maximize the classification loss

Figure 4: The illustration of the nodule synthesis framework.

network  $f_{FPR}$  outputs the classification evidence  $e_i$  for positive and negative labels.  $L_{FPR}(e_i, 1)$  is the classification loss defined with the beta distribution distance. For the symbolic brevity, we refer to [29] for the detailed definition of the beta distribution network output and the loss function. It was shown in [29] that the beta distribution networks are less likely to be activated by out-of-distribution (OOD) patches and can produce comparable classification accuracy as cross-entropy. However, in our experiments, we also show that uncertainty estimation alone does not suffice to make the network robust to all OOD samples especially when such samples are searched by PGD. The perturbation vector  $\delta$  can be updated as

$$\delta := \mathcal{P}(\delta + \alpha \nabla_{\delta} L_{FPR}(\cdot)) \quad (9)$$

where  $\mathcal{P}$  denotes the projection onto the ball of interest defined by  $\epsilon$ ;  $\alpha$  is the step size. In our experiments, we set  $\epsilon = 0.15$ ,  $\alpha = 0.05$ .  $\delta$  is updated with 20 iterations for each search. In Fig. 5, we show initial synthetic nodules together with the synthetic nodules searched with PGD. Though visually similar, the tiny differences in the nodule appearance can result in a large difference in the  $F_{FPR}$  responses.

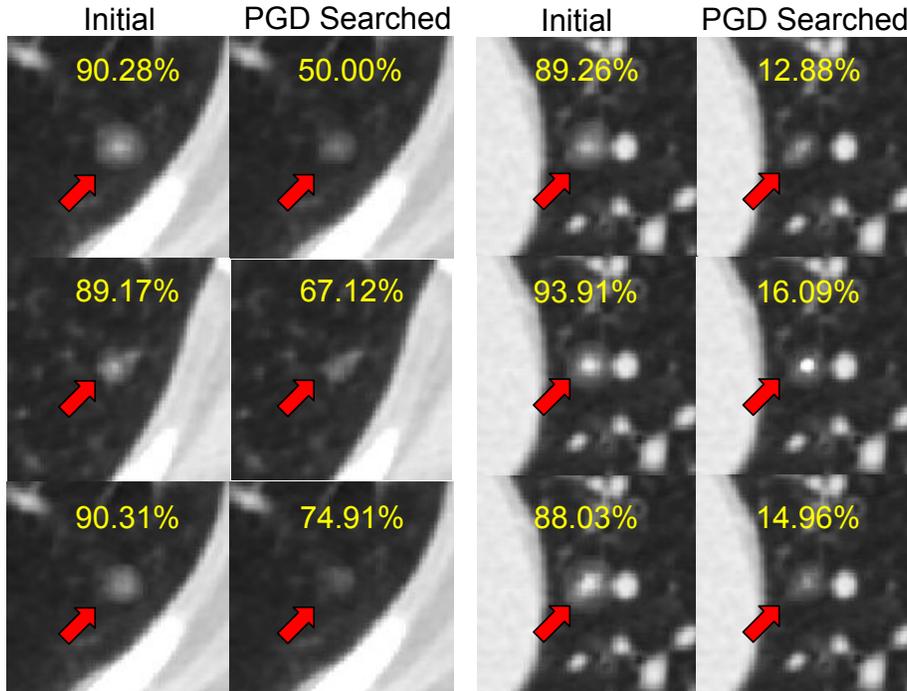


Figure 5: The demonstrations of the synthetic nodules before and after PGD searching. With slight perturbation in the nodule appearance, the nodule detector trained with conventional strategy would output significantly lower confidence score.

### 3.3 Over-confident Perturbation with PGD Sampling

Besides searching the latent codes for the nodule synthesizer, PGD can also be used for perturbing the real patches  $x_i$  as

$$e_i = f_{FPR}(x_i + \delta_i^{patch}) \quad (10)$$

$$\arg \max_{\|\delta_i^{patch}\| \leq \epsilon} L_{FPR}(e_i, g_i) \quad (11)$$

where  $g_i$  is the groundtruth label for patch  $x_i$ . As shown in the first row of Fig. 6, we found for most of the positive nodules patches, it is easy to find a  $\delta_i^{patch}$  with a small magnitude to perturb  $x_i$  so that  $f_{FPR}$  no longer recognizes the nodule resides in it. Such perturbations can disturb the model from recognizing the nodules when the images contain unexpected abnormalities, strong imaging artefacts or malicious noise injections.

We also found that even for noise patches  $x_i^{uniform}$  drawn from a uniform distribution, PGD can search for a neighbouring patch and excites the FPR network to output a positive decision, though the searched patch does not contain any interpretable patterns as shown in the second row of Fig. 6. The intersection between the chest CT distribution and the uniform distribution is expected to have close to zero probability mass. As

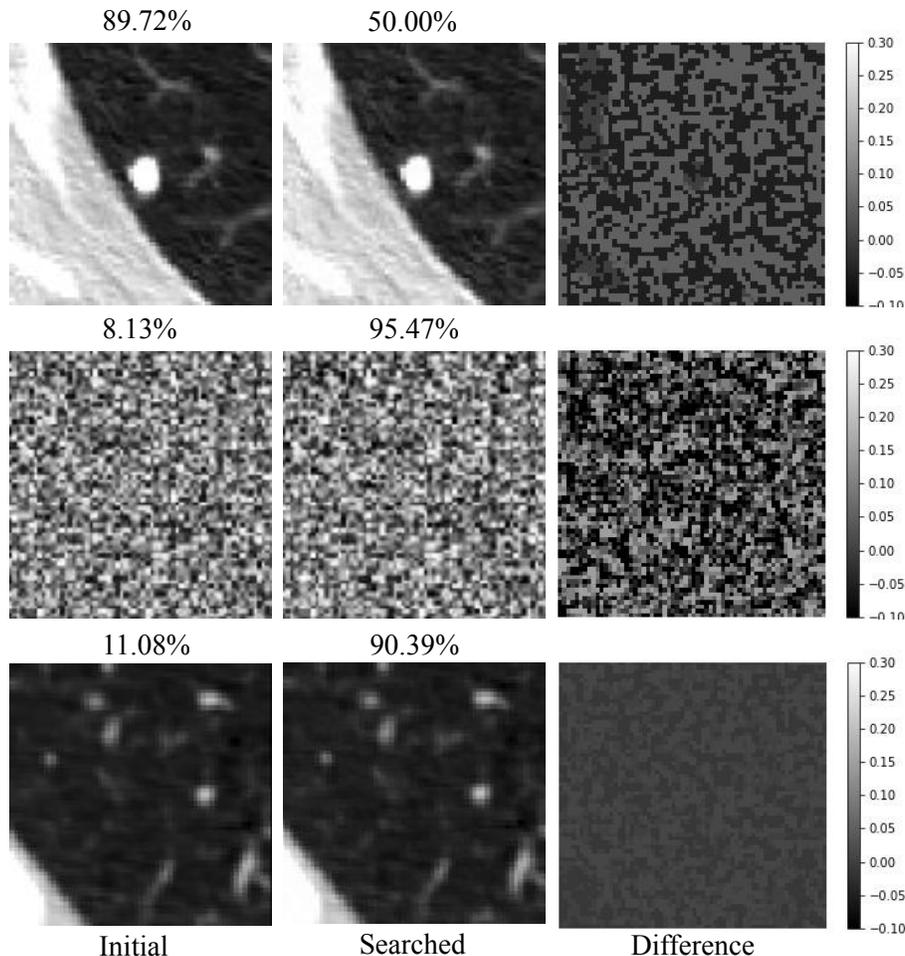


Figure 6: The upper row demonstrates the noise perturbation on nodule patches. Arbitrary noise can determine a trained nodule detector to ignore a well-defined nodule. The middle and bottom row demonstrates that specific noise patterns can activate a trained nodule detector to output high confidence scores from either pure adversarial noise or the negative CT patches distorted by adversarial noise. The difference patches are shown with the window  $[0, 0.3]$  to make the perturbation visible while the image patches are shown with the window  $[0, 1]$ .

explained in [5], ReLU networks decompose the observation space into a finite set of polytopes in which outer polytopes extend to infinity. Adding the adversarial patches searched by Eq.(11) to augment the FPR network can make it robust to such image perturbations by closing the decision boundary.

In practice, we train a baseline FPR network first by randomly sampling real positive and negative candidate patches with 50% chance each until reaching convergence.

Then we finetune the baseline model by also sampling from the augmentation patches generated by attacking the baseline model. For positive sampling, we draw 50% from the real positive patches, and 25% from synthetic nodules and 25% the adversarial positive patches. For negative sampling, we draw 50% from both real negative patches and 50% from the adversarial negative patches.

## 4 Data and Experiment Settings

6488 3D chest CT scans were collected for training. The training images were collected from multiple sources, including the LUNA challenge [6], the NLST cohort [14] and an in-house data collection. Each training image contains at least one radiologist confirmed nodule. We annotated the nodule locations and diameters in the training images from our in-house dataset and the NLST subset. Our annotators firstly detected all the potential nodule candidates. Then two radiologists went through all the candidates to confirm the presence of a nodule. 10% of the training images were randomly sampled as the validation set for parameter searching and early stopping. To evaluate the performance, we constructed two benchmark datasets, as summarized in Table 1. The In-house Benchmark was built based on a private data collection with 174 challenging CT scans. Besides lung nodules, many patients in the In-house Benchmark also had other types of pulmonary abnormalities which constitute a significant source of false positives for both human readers and the networks. The NLST Benchmark consists of randomly sampled 272 baseline CT scans from the NLST cohort. The patients were sampled following the real-world screening distribution [14] (1% with cancer, 25.8% with cancer negative nodules and 73.2% healthy) while ensuring (1) the slice thicknesses are lower than  $1.5mm$  (2) there is no gap in the DICOM series (3) each image contains the entire lung. We had three on-board radiologists read the images in both benchmark datasets independently. In the first round, each radiologist marked the nodule candidates individually. All the candidate nodules spotted in the first round were merged and presented to each radiologist to confirm in case there were under-attended nodule candidates. We took the nodules that are the consensus among all three radiologists as the positive locations while the rest as irrelevant findings which were not involved in the metrics computing. We only considered the nodules with the diameters larger than  $6mm$  for benchmarking. However, we do not claim this is a critical choice since the size threshold can be adjusted according to the different application scenarios.

All the augmentation patches, including the synthetic nodules, perturbed positive nodule patches and the perturbation noises, were pre-computed and randomly sampled during the FPR model training by attacking the baseline network (baseline-beta-finetune). Therefore, the large pool of randomly generated patches was kept consistent across different experiments to guarantee the reproducibility of all the experiment results. We generated synthetic nodule patches on 10 random background patches from each training image. The locations of the background patches were constrained within the lungs using the lung segmentation masks predicted by a previously trained network. We also ensured that the background patches do not contain a real nodule inside. For each background patch, we sampled the synthesizer six times with random sampling and the PGD sampling, respectively. It resulted in 389,280 synthetic nodule patches

Table 1: The summary of the chest CT benchmark datasets.

	In-house Benchmark	NLST Benchmark
CT scans	174	272
CT scans w/ Nodules	97	83
Solid Nodules	94	103
Fully Calcified Nodules	7	19
Part-Solid Nodules	13	3
Ground Glass Nodules	36	6
Total Nodules ( $\geq 6\text{mm}$ )	150	131

for both sampling strategies. We generated one adversarially perturbed patch for each positive nodule candidate in our training data (22,169 relevant nodules) similarly to the upper row of Fig. 6. We also generated 100,000 pure adversarial noise patches similarly to the lower row of Fig. 6. To stress-test the robustness of network at random pulmonary locations, we sampled 10 random patches centered in the lungs as the negative stress-test samples from each benchmark CT volume, while avoiding annotated nodules. We add adversarial noise to these negative samples by attacking the baseline network (baseline-beta-finetune). For all the experiments involving Poisson noise, we used the same ratio to sample the Poisson noise injected patches as used for the adversarial noise patches.

## 5 Results

### 5.1 Toy Example

In Fig. 7, we firstly show a toy experiment built with the simple two-moon dataset to demonstrate the presented concept. 500 spots are sampled from both the positive and the negative cluster by adding the Gaussian noise with the standard deviation of 0.15. In our context, they represent the positive and negative candidates used for training the FPR classifier. We train a ReLU activated multi-layer perceptron to mimic the FPR classifier based on the sampled spots to plot the decision boundary. We then sub-sample only 20 positive candidates following a long tail distribution to simulate the real-world training set distribution as Fig. 7b. We trained a small VAE on the 20 positive spots and generated synthetic samples by drawing the latent code from a standard normal distribution. The added synthetic spots help filling the hole in the decision boundary as in Fig. 7c. However, a sizeable out-of-distribution area is also predicted as confident positive as anticipated in [5]. We then sampled another 20 spots that are randomly drawn from a uniform distribution and added them to the negative cluster. In Fig. 7d, it is shown that such noise samples can bound the decision boundary tightly to the positive cluster. Though there is a small chance that the noise spots can also reside in the positive cluster, such cases are extremely rare in the real world 3D inputs. Though we use uniform sampling in this toy example, it is notable that in a high-dimensional input space, the random sampling can be highly in-efficient for both

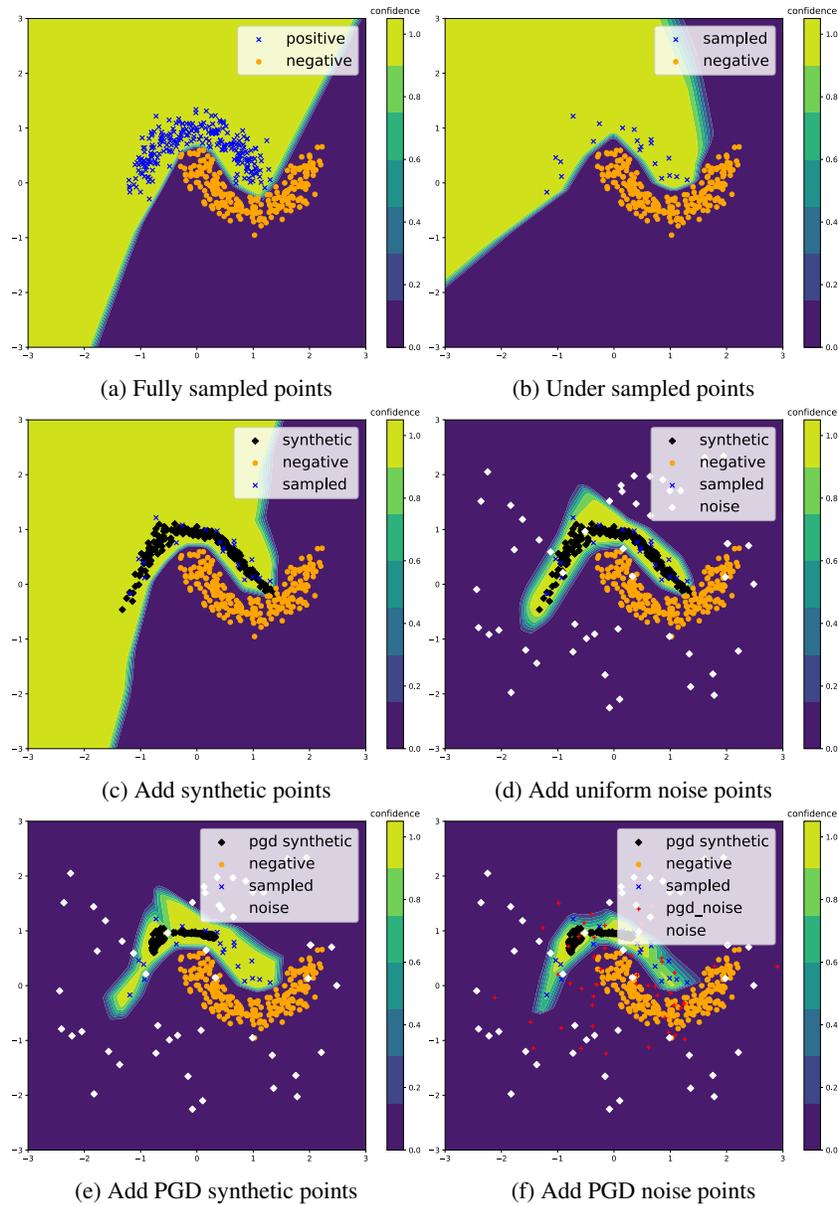


Figure 7: A toy experiment to depict the concept of the proposed augmentation methods.

synthesizing real nodules and generating adversarial noise samples. We use PGD to search for the latent code from the trained VAE. As in Fig. 7e, the PGD searched synthetic spots only reside in the under-sampled region. In addition to the uniform

Table 2: The table summarizes the complexities of the nodule detection models.

Network	#Param	Mac (G)	Input Size
DeepLung [10]	5.36M	168.43	128×128×128
Blob Solid / GGO ResUNet	143,92k	67.15	128×128×128
RPN	276.93k	113.91	129×128×128
FPR	455.98k	17.24	64×64×64

Table 3: The table summarizes the FROC metrics on the clean benchmark datasets. The CPM score averages the sensitivities sampled at 7 log-scale operating points indicating different numbers of false positives (0.125, 0.25, 0.5, 1, 2, 4, 8).

In-house Benchmark											
	PERTURB	SYN	LOSS	CPM	FP=0.125	FP=0.25	FP=0.5	FP=1	FP=2	FP=4	FP=8
baseline-DeepLung [10]	✗	✗	N/A	55.15% (46.56% - 63.92%)	31.59%	38.81%	51.84%	58.64%	62.98%	68.59%	73.58%
baseline-ce	✗	✗	CE	88.46% (83.18% - 93.76%)	73.09%	82.24%	89.84%	92.02%	92.28%	94.65%	96.64%
baseline-beta	✗	✗	BETA	89.11% (82.91% - 93.87%)	75.43%	83.93%	88.78%	90.80%	91.77%	94.26%	96.64%
baseline-beta-finetune	✗	✗	BETA	88.90% (83.62% - 94.22%)	75.58%	84.29%	90.79%	91.37%	92.11%	94.24%	96.70%
beta+syn (random)	✗	✓	BETA	90.76% (85.95% - 95.49%)	79.89%	87.43%	92.05%	93.34%	93.35%	94.26%	96.67%
beta+syn	✗	✓	BETA	91.22% (86.13% - 95.69%)	81.09%	87.66%	<b>92.66%</b>	93.33%	93.33%	94.17%	96.99%
beta+poisson	✗	✓	BETA	91.98% (86.74% - 96.30%)	81.50%	87.24%	90.62%	93.40%	<b>95.87%</b>	<b>97.29%</b>	97.96%
beta+poisson+syn	✗	✓	BETA	<b>92.26%</b> ( <b>87.52%</b> - <b>96.30%</b> )	<b>81.96%</b>	<b>88.16%</b>	92.65%	<b>94.00%</b>	95.10%	96.04%	<b>98.00%</b>
beta+perturb	✓	✗	BETA	90.07% (85.19% - 95.02%)	76.35%	85.83%	89.91%	93.42%	93.61%	95.40%	97.92%
beta+perturb+syn	✓	✓	BETA	90.47% (85.22 - 95.12)	77.52%	86.29%	89.95%	92.75%	93.97%	94.98%	97.45%
NLST Benchmark											
	PERTURB	SYN	LOSS	CPM	FP=0.125	FP=0.25	FP=0.5	FP=1	FP=2	FP=4	FP=8
baseline-DeepLung [10]	✗	✗	N/A	75.71% (67.16% - 84.32%)	32.77%	45.05%	52.14%	57.07%	67.21%	74.09%	75.71%
baseline-ce	✗	✗	CE	82.56% (73.26% - 91.03%)	52.18%	71.50%	84.99%	89.68%	91.68%	93.14%	95.63%
baseline-beta	✗	✗	BETA	80.62% (69.49% - 89.78%)	44.74%	68.39%	83.35%	88.56%	91.38%	93.18%	94.40%
baseline-beta-finetune	✗	✗	BETA	83.60% (74.19% - 91.91%)	53.69%	73.44%	85.30%	91.35%	93.01%	93.99%	<b>95.71%</b>
beta+syn (random)	✗	✓	BETA	85.81% (77.40% - 93.43%)	66.04%	80.10%	86.55%	90.17%	92.05%	93.15%	95.06%
beta+syn	✗	✓	BETA	<b>87.89%</b> ( <b>81.20%</b> - <b>93.93%</b> )	<b>74.44%</b>	<b>81.04%</b>	87.61%	<b>90.55%</b>	93.30%	93.80%	94.77%
beta+poisson	✗	✓	BETA	86.33% (76.71% - 93.46%)	64.95%	78.22%	87.41%	90.25%	<b>93.85%</b>	<b>94.18%</b>	95.44%
beta+poisson+syn	✗	✓	BETA	86.55% (78.16% - 93.52%)	68.25%	78.33%	86.96%	90.46%	93.14%	93.84%	94.89%
beta+perturb	✓	✗	BETA	85.63% (76.91% - 93.43%)	64.92%	79.83%	86.50%	89.23%	92.63%	93.85%	94.67%
beta+perturb+syn	✓	✓	BETA	86.71% (76.44% - 93.74%)	64.47%	80.16%	<b>88.00%</b>	89.96%	92.62%	93.77%	94.52%

spots, we show the PGD searched negative spots which are closer to the positive cluster in Fig. 7f. Such supporting negative spots can be more efficient for refining the decision boundary when the input dimension is higher as in 3D chest CT patches.

## 5.2 Benchmark on clean data

Before we analyze the FPR networks, the frozen CG framework achieved 100% sensitivity on the In-house Benchmark and 97.71% sensitivity on the NLST Benchmark when having 100 average candidates per scan. We summarize the FROC curves for benchmarking the nodule detection FPR models trained with different strategies in Table 3. Similar LUNA16 [6], the final CPM score is defined as the average sensitivity at 7 log-scaled false positive rates: 0.125, 0.25, 0.5, 1, 2, 4, and 8 FPs per scan. The performance obtained by training the DeepLung<sup>1</sup> [10] (55.15% CPM and 75.71% CPM) is used as a reference standard for the detection performance for the baseline detection frameworks on our benchmark datasets. Notably, our In-house benchmark has larger

<sup>1</sup><https://github.com/wentaozhu/DeepLung>

variances in nodule types, sizes and contexts than the LIDC dataset [13]. The classification head with beta distribution (baseline-beta) produced similar CPM scores as the sigmoid head trained with binary cross-entropy (baseline-ce). However, we also show that the classifier would generate slightly higher CPM scores if the network is firstly trained with cross-entropy and then finetuned with the beta-distribution loss (baseline-beta-finetune). In the experiments beta-syn (random) and beta+syn, we respectively added synthetic nodules randomly sampled from the standard normal, and the ones searched using the proposed PGD sampling. Though both types of synthetic nodules can improve the overall network generalization, the nodules searched with PGD consistently outperforms its counterpart, especially at the region of the lower number of false positives. We show that adding the noise perturbation augmentation patches (beta+perturb and beta+perturb+syn) can also slightly improve the overall CPM scores comparing to the conventional training baseline (baseline-beta-finetune). However, they do not show better performance than only using only PGD searched nodules (beta+syn). We show such perturbation augmented networks are more robust to both uniform and adversarial noise in the next section. The data augmentation by adding random Poisson noise [41] into the training patches (beta+poisson and beta+poisson+syn) achieved slightly higher detection performance than adding the adversarial noise in both clean benchmarks. The complexities of the detection models are summarized in Table 2.

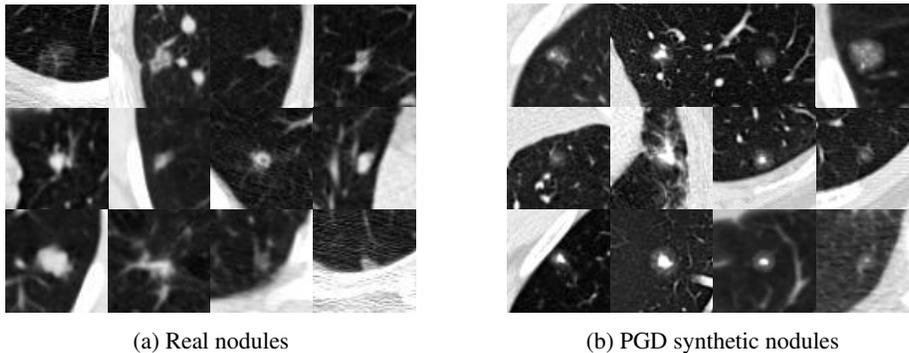


Figure 8: The mosaic view to compare the real nodule patches and the synthetic nodules in patches of size  $64^3$  and  $0.625^3$ mm resolution. Besides being generally smaller, the PGD searched synthetic nodules tend to have round glass component with or without a solid core. Such non-solid or part-solid nodules are relatively rare in the real datasets.

### 5.3 Stress test

#### 5.3.1 Synthetic nodules

The central slices of the randomly selected real nodules and the hard nodules sampled by PGD are shown in Fig. 8. Besides being generally smaller, the PGD searched synthetic nodules tend to have round glass component with or without a solid core. Such non-solid or part-solid nodules are relatively rare in the real datasets. Though one can

Table 4: False positive reduction confidence means and standard deviations obtained on the synthetic nodule generated by using randomly sampled coding (Randomly sampled nodules) and the PGD sampled coding (Randomly sampled nodules).

Method	Randomly sampled nodules	PGD sampled nodules
baseline-beta-finetune	$0.82 \pm 0.17$	$0.53 \pm 0.15$
beta+syn	<b><math>0.92 \pm 0.09</math></b>	<b><math>0.88 \pm 0.16</math></b>

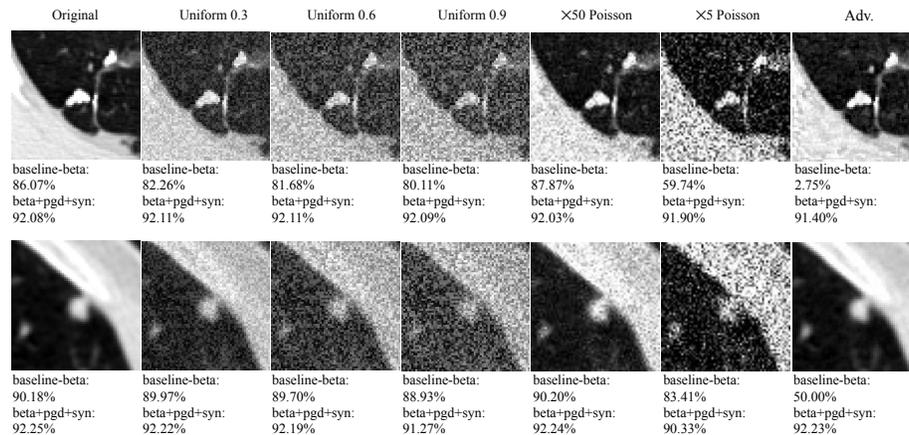


Figure 9: Examples to show different levels of uniform noise and adversarial noise on two nodules randomly drawn from the stress-test.

still visually distinguish a subset of the synthetic nodules from the real nodules, they can be a valuable source to stress-test the FPR network as most of such cases reside at the original decision boundaries. We synthesized 10000 nodules with both random Gaussian sampling and PGD searching respectively. They were fed to the FPR networks trained with (beta+syn) and without (baseline-beta-finetune) synthetic nodules. We ensured that all the synthetic nodules in this test have diameters at least 6mm. The mean and standard deviations of the network responses are shown in Table. 4. Though the conventional network achieved 88.90% CPM, it failed to recognize many nodules generated with only random sampling synthesis. The baseline network predicts the majority of the PGD searched nodules around 0.53, which are defined as out-of-distribution (OOD) samples. The network augmented with PGD synthetic nodules can successfully recognize most of the PGD synthetic nodules with high-confidence with mean confidence 0.87.

### 5.3.2 Noise

To stress-test the network resistance to different levels of noise, we first add uniform noise with different magnitudes to the nodule patches as depicted by Fig. 9. The uniform noise can reduce the baseline network response from 0.86 to 0.81 as shown in Ta-

Table 5: Stress-test by perturbing the positive patches with different levels of uniform, Poisson noise and PGD noise perturbation.

Method	0.0 uniform noise	0.3 uniform noise	0.6 uniform noise	0.9 uniform noise	$\times 50$ Poisson noise	$\times 5$ Poisson noise	Adv.
baseline-beta-finetune	0.86 $\pm$ 0.12	0.86 $\pm$ 0.11	0.83 $\pm$ 0.14	0.81 $\pm$ 0.16	0.86 $\pm$ 0.11	0.78 $\pm$ 0.18	0.25 $\pm$ 0.25
beta+syn	0.91 $\pm$ 0.10	0.91 $\pm$ 0.10	0.89 $\pm$ 0.11	0.87 $\pm$ 0.16	0.91 $\pm$ 0.09	0.84 $\pm$ 0.19	0.37 $\pm$ 0.39
beta+perturb	0.91 $\pm$ 0.06	0.91 $\pm$ 0.05	<b>0.91 <math>\pm</math> 0.05</b>	<b>0.89 <math>\pm</math> 0.12</b>	0.91 $\pm$ 0.05	<b>0.91 <math>\pm</math> 0.07</b>	<b>0.91 <math>\pm</math> 0.07</b>
beta+perturb+syn	<b>0.92 <math>\pm</math> 0.04</b>	<b>0.91 <math>\pm</math> 0.04</b>	0.90 $\pm$ 0.07	0.87 $\pm$ 0.13	<b>0.92 <math>\pm</math> 0.04</b>	0.90 $\pm$ 0.07	0.90 $\pm$ 0.08
beta+poisson	0.86 $\pm$ 0.07	0.86 $\pm$ 0.07	0.84 $\pm$ 0.09	0.83 $\pm$ 0.11	0.86 $\pm$ 0.07	0.80 $\pm$ 0.14	0.78 $\pm$ 0.18
beta+poisson+syn	0.90 $\pm$ 0.07	0.89 $\pm$ 0.09	0.88 $\pm$ 0.11	0.86 $\pm$ 0.14	0.89 $\pm$ 0.07	0.79 $\pm$ 0.21	0.82 $\pm$ 0.17

Table 6: Stress test by feeding noise patches and the negative patches to the network.

Method	Uniform noise	PGD noise	Negative patches	PGD negative patches
baseline-beta-finetune	0.11 $\pm$ 0.05	0.87 $\pm$ 0.17	0.14 $\pm$ 0.09	0.82 $\pm$ 0.21
beta+syn	0.22 $\pm$ 0.16	0.90 $\pm$ 0.10	<b>0.10 <math>\pm</math> 0.11</b>	0.63 $\pm$ 0.32
beta+perturb	<b>0.04 <math>\pm</math> 0.00</b>	<b>0.08 <math>\pm</math> 0.13</b>	0.19 $\pm$ 0.14	0.26 $\pm$ 0.20
beta+perturb+syn	0.05 $\pm$ 0.01	0.09 $\pm$ 0.09	0.21 $\pm$ 0.15	0.32 $\pm$ 0.23
beta+poisson	0.13 $\pm$ 0.06	0.23 $\pm$ 0.21	0.13 $\pm$ 0.12	0.18 $\pm$ 0.18
beta+poisson+syn	0.08 $\pm$ 0.02	0.34 $\pm$ 0.25	0.11 $\pm$ 0.10	<b>0.17 <math>\pm</math> 0.17</b>

ble. 5. We found that the network augmented with either synthetic nodules (beta+syn) or PGD noises (beta+perturb and beta+perturb+syn) can be more robust to uniform noise. To simulate the Poisson noise in CT, we rescaled the CT patches to  $[0, 50]$  and  $[0, 1]$  respectively and then sample from them following the Poisson process. Similarly to the uniform noise, stronger ( $\times 5$ ) Poisson noise can deactivate the baseline FPR network (from 0.86 to 0.78) while affecting less on the augmented networks augmented with PGD noise. Interestingly, the networks augmented with Poisson noise (beta+poisson and beta+poisson+syn) did not show much better robustness towards Poisson noise than the baseline network. The networks trained without adversarial noise augmentation tend to be deactivated by the adversarial noise. At the same time, the mean responses from the adversarial noise augmented networks remained around 0.9. We also tested the FPR network by feeding randomly generated noise patches and PGD adversarial noise patches as shown in Table. 6. The conventional FPR network would normally not be activated by random uniform noise, meaning most of the mean responses are below 0.22. However, pure adversarial noise patches can easily activate the baseline networks (0.87 and 0.90). The networks augmented with the Poisson noise show stronger robustness than the baseline networks (0.23 and 0.34). The networks augmented by the adversarial noise augmentation (beta+perturb and beta+perturb+syn) are more robust to both types of noise patterns with below 0.1 mean responses. In Table. 6 we also show that the networks augmented with either the adversarial noise and Poisson noise are more robust to the adversarial noise added to the real negative CT patches than the baseline network.

## 6 Discussions and Conclusions

In this paper, we propose adversarial augmentation methods to improve the robustness of the nodule detection framework against two major sources of the out-of-distribution

samples (1) the nodules with under-represented properties in the training dataset (2) the images with unexpected noise or contrast. We first use the beta-distribution to replace the sigmoid output of the false positive reduction network to estimate the observation uncertainty explicitly at the output layer. Then we add both adversarial synthetic nodules and adversarial perturbation noise to the training set that is searched using the project gradient descent (PGD). Some of the existing works as listed in S.II.B have attempted to use synthetic images to improve the model performance since generative models can provide samples under-represented in the training set which is not easily simulated by the conventional data augmentation techniques. However, as shown in Table II, the performance gain using randomly sampled synthetic images is limited since the majority of the randomly sampled cases drawn from a prior distribution are well represented by the training data distribution. Therefore, we propose to search for the latent code associated with the hard-samples from the synthesizer using PGD to maximize the training loss of the supervised model. This applies when (1) the data synthesizer and the supervised learning model are both differentiable (2) the data synthesizer is trained to well-represent the manifold of interests. We show that the proposed techniques can improve the generalization of the nodule detector by learning pathologically relevant patterns, we tested it on two benchmark datasets with groundtruth annotated by experienced radiologists. We also use the synthetic nodules and the generated perturbations to stress test the trained models and show the augmented networks can be more robust to both hard nodules as well as different types of noise distortions. By using the beta distribution based uncertainty estimation, we also showed that uncertainty estimation alone might not be sufficient to make the network robust to the out-of-distribution inputs, especially when the inputs are adversarially generated. The proposed framework can be indeed applied to similar classification, segmentation or detection models as long as the model and the data synthesizer are fully differentiable.

Some studies observed that adversarial samples can cause decreased testing performance on clean testing data [33, 34, 35, 36]. However, bearing in mind that dropped accuracy can be expected as a cost for robustness, in our experiments shown in Table. 3, we did not notice obvious testing performance decreases on the clean data though the best performance is indeed observed on the models with only adversarial hard-case sampling without adversarial noise. This observation is also consistent between the two benchmark datasets. Our assumptions are two-fold: (1) As hypothesized in [36], the generalization performance drop might be because the clean data and the adversarial data are drawn from different manifolds. Even though detecting rare nodules automatically has been a challenging problem for decades, the underlying data manifolds of the nodules can be much simpler than the ones of natural image applications. The neural networks can, therefore, be over-parameterized to fit both the clean and the adversarial manifold without changing the normalization layers. (2) Other than adversarial noise, the adversarial noise enhanced networks can be also more robust to other noise or artifact sources that can appear in clean CT data. However, these assumptions might need to be further validated by future work.

As one of the early attempts to enhance the robustness of the medical image analysis CNNs, this study has a few limitations that would be targeted in future works. We use a relatively simple nodule synthesizer network and the standard PGD to sample the lung nodules from the latent space. This synthesizer was not capable of synthesizing all

types of different nodules, such as nodules with spiculation. It was also not constrained to maintain the size of a synthetic nodule, therefore we had to filter out the synthetic nodules that are smaller than the relevant threshold. We only investigated the network robustness towards three types of image noise. The improved robustness towards other types of image artefacts, such as metal artefacts and motion distortion, etc., remains unknown. As a proof of concept study, the proposed techniques were only applied to the false positive reduction (FPR) of the lung nodule detection pipeline for brevity. However, the same perturbations can also affect candidate generation networks. We also found that in practice it is hard to generate adversarial noise by attacking the noise augmented networks without showing visually detectable artefacts. However, it is possible to attack augmented networks with the same techniques. Though we only evaluated the proposed techniques in the context of nodule detection, we believe such techniques can also be helpful for the other deep CNN based medical imaging applications with minor technical adjustments.

**Disclaimer:** The concepts and information presented in this paper are based on research results that are not commercially available

## References

- [1] R. L. Siegel, K. D. Miller, and A. Jemal, “Cancer statistics, 2019,” *CA: a cancer journal for clinicians*, vol. 69, no. 1, pp. 7–34, 2019.
- [2] D. M. Hansell, A. A. Bankier, H. MacMahon, T. C. McLoud, N. L. Muller, and J. Remy, “ Fleischner society: glossary of terms for thoracic imaging,” *Radiology*, vol. 246, no. 3, pp. 697–722, 2008.
- [3] R. J. McCunney and J. Li, “Radiation risks in lung cancer screening programs,” *Chest*, vol. 145, no. 3, pp. 618–624, 2014.
- [4] H. J. de Koning, C. M. van der Aalst, P. A. de Jong, E. T. Scholten, K. Nackaerts, M. A. Heuvelmans, J.-W. J. Lammers, C. Weenink, U. Yousaf-Khan, N. Horeweg, S. van ’t Westeinde, M. Prokop, W. P. Mali, F. A. Mohamed Hoesein, P. M. van Ooijen, J. G. Aerts, M. A. den Bakker, E. Thunnissen, J. Verschakelen, R. Vliegenthart, J. E. Walter, K. ten Haaf, H. J. Groen, and M. Oudkerk, “Reduced lung-cancer mortality with volume ct screening in a randomized trial,” *New England Journal of Medicine*, vol. 382, no. 6, pp. 503–513, 2020, pMID: 31995683.
- [5] M. Hein, M. Andriushchenko, and J. Bitterwolf, “Why ReLU networks yield high-confidence predictions far away from the training data and how to mitigate the problem,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 41–50.
- [6] A. A. A. Setio, A. Traverso, T. De Bel, M. S. Berens, C. van den Bogaard, P. Cerello, H. Chen, Q. Dou, M. E. Fantacci, B. Geurts *et al.*, “Validation, comparison, and combination of algorithms for automatic detection of pulmonary nod-

- ules in computed tomography images: the LUNA16 challenge,” *Medical image analysis*, vol. 42, pp. 1–13, 2017.
- [7] A. A. A. Setio, F. Ciompi, G. Litjens, P. Gerke, C. Jacobs, S. J. van Riel, M. M. W. Wille, M. Naqibullah, C. I. Sánchez, and B. van Ginneken, “Pulmonary nodule detection in CT images: false positive reduction using multi-view convolutional networks,” *IEEE transactions on medical imaging*, vol. 35, no. 5, pp. 1160–1169, 2016.
- [8] G. Pérez and P. Arbeláez, “Automated detection of lung nodules with three-dimensional convolutional neural networks,” in *13th International Conference on Medical Information Processing and Analysis*, vol. 10572. International Society for Optics and Photonics, 2017, p. 1057218.
- [9] Q. Dou, H. Chen, L. Yu, J. Qin, and P.-A. Heng, “Multilevel contextual 3-D CNNs for false positive reduction in pulmonary nodule detection,” *IEEE Transactions on Biomedical Engineering*, vol. 64, no. 7, pp. 1558–1567, 2016.
- [10] W. Zhu, C. Liu, W. Fan, and X. Xie, “DeepLung: Deep 3D dual path nets for automated pulmonary nodule detection and classification,” in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2018, pp. 673–681.
- [11] H. Xie, D. Yang, N. Sun, Z. Chen, and Y. Zhang, “Automated pulmonary nodule detection in CT images using deep convolutional neural networks,” *Pattern Recognition*, vol. 85, pp. 109–119, 2019.
- [12] J. Ding, A. Li, Z. Hu, and L. Wang, “Accurate pulmonary nodule detection in computed tomography images using deep convolutional neural networks,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2017, pp. 559–567.
- [13] S. G. Armato et al., “The Lung Image Database Consortium (LIDC) and Image Database Resource Initiative (IDRI): A completed reference database of lung nodules on CT scans,” *Medical Physics*, vol. 38, no. 2, pp. 915–931, 1 2011.
- [14] N. L. S. T. R. Team, “The national lung screening trial: overview and study design,” *Radiology*, vol. 258, no. 1, pp. 243–253, 2011.
- [15] D. Heaven, “Why deep-learning ais are so easy to fool,” *Nature*, vol. 574, no. 7777, p. 163, 2019.
- [16] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, “Towards deep learning models resistant to adversarial attacks,” *arXiv preprint arXiv:1706.06083*, 2017.
- [17] J. Zhang, Y. Xia, H. Cui, and Y. Zhang, “Pulmonary nodule detection in medical images: a survey,” *Biomedical Signal Processing and Control*, vol. 43, pp. 138–147, 2018.

- [18] L. M. Pehrson, M. B. Nielsen, and C. Ammitzbøl Lauridsen, “Automatic pulmonary nodule detection applying deep learning or machine learning algorithms to the LIDC-IDRI database: A systematic review,” *Diagnostics*, vol. 9, no. 1, p. 29, 2019.
- [19] X. Yi, E. Walia, and P. Babyn, “Generative adversarial network in medical imaging: A review,” *Medical image analysis*, p. 101552, 2019.
- [20] J. Yang, S. Liu, S. Grbic, A. A. A. Setio, Z. Xu, E. Gibson, G. Chabin, B. Georgescu, A. F. Laine, and D. Comaniciu, “Class-aware adversarial lung nodule synthesis in CT images,” in *2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019)*. IEEE, 2019, pp. 1348–1352.
- [21] S. Liu, E. Gibson, S. Grbic, Z. Xu, A. A. A. Setio, J. Yang, B. Georgescu, and D. Comaniciu, “Decompose to manipulate: Manipulable object synthesis in 3D medical images with structured image decomposition,” *arXiv preprint arXiv:1812.01737*, 2018.
- [22] D. Jin, Z. Xu, Y. Tang, A. P. Harrison, and D. J. Mollura, “CT-realistic lung nodule simulation from 3D conditional generative adversarial networks for robust lung segmentation,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2018, pp. 732–740.
- [23] Z. Xu, X. Wang, H.-C. Shin, H. Roth, D. Yang, F. Milletari, L. Zhang, and D. Xu, “Tunable CT lung nodule synthesis conditioned on background image and semantic features,” in *International Workshop on Simulation and Synthesis in Medical Imaging*. Springer, 2019, pp. 62–70.
- [24] Z. Xu, X. Wang, H.-C. Shin, D. Yang, H. Roth, F. Milletari, L. Zhang, and D. Xu, “Correlation via synthesis: end-to-end nodule image generation and radiogenomic map learning based on generative adversarial network,” *arXiv preprint arXiv:1907.03728*, 2019.
- [25] C. Gao, S. Clark, J. Furst, and D. Raicu, “Augmenting LIDC dataset using 3D generative adversarial networks to improve lung nodule detection,” in *Medical Imaging 2019: Computer-Aided Diagnosis*, vol. 10950. International Society for Optics and Photonics, 2019, p. 109501K.
- [26] C. Han, Y. Kitamura, A. Kudo, A. Ichinose, L. Rundo, Y. Furukawa, K. Umemoto, Y. Li, and H. Nakayama, “Synthesizing diverse lung nodules wherever massively: 3D multi-conditional GAN-based CT image augmentation for object detection,” in *2019 International Conference on 3D Vision (3DV)*. IEEE, 2019, pp. 729–737.
- [27] Q. Wang, X. Zhou, C. Wang, Z. Liu, J. Huang, Y. Zhou, C. Li, H. Zhuang, and J.-Z. Cheng, “WGAN-based synthetic minority over-sampling technique: improving semantic fine-grained classification for lung nodules in CT images,” *IEEE Access*, vol. 7, pp. 18 450–18 463, 2019.

- [28] C. Mayer and R. Timofte, “Adversarial sampling for active learning,” in *The IEEE Winter Conference on Applications of Computer Vision*, 2020, pp. 3071–3079.
- [29] F. C. Ghesu, B. Georgescu, E. Gibson, S. Guendel, M. K. Kalra, R. Singh, S. R. Digumarthy, S. Grbic, and D. Comaniciu, “Quantifying and leveraging classification uncertainty for chest radiograph assessment,” in *Medical Image Computing and Computer Assisted Intervention – MICCAI 2019*, D. Shen, T. Liu, T. M. Peters, L. H. Staib, C. Essert, S. Zhou, P.-T. Yap, and A. Khan, Eds. Cham: Springer International Publishing, 2019, pp. 676–684.
- [30] M. Sensoy, L. Kaplan, and M. Kandemir, “Evidential deep learning to quantify classification uncertainty,” in *Advances in Neural Information Processing Systems 31*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds. Curran Associates, Inc., 2018, pp. 3179–3189.
- [31] A. Galloway, A. Golubeva, T. Tanay, M. Moussa, and G. W. Taylor, “Batch normalization is a cause of adversarial vulnerability,” *arXiv preprint arXiv:1905.02161*, 2019.
- [32] A. Meinke and M. Hein, “Towards neural networks that provably know when they don’t know,” *arXiv preprint arXiv:1909.12180*, 2019.
- [33] D. Stutz, M. Hein, and B. Schiele, “Disentangling adversarial robustness and generalization,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6976–6987.
- [34] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry, “Robustness May Be at Odds with Accuracy,” *arXiv e-prints*, p. arXiv:1805.12152, May 2018.
- [35] A. Raghunathan, S. M. Xie, F. Yang, J. C. Duchi, and P. Liang, “Adversarial training can hurt generalization,” *arXiv preprint arXiv:1906.06032*, 2019.
- [36] C. Xie and A. Yuille, “Intriguing properties of adversarial training at scale,” in *International Conference on Learning Representations*, 2020.
- [37] M. Paschali, S. Conjeti, F. Navarro, and N. Navab, “Generalizability vs. Robustness: Adversarial Examples for Medical Imaging,” *arXiv e-prints*, p. arXiv:1804.00504, Mar 2018.
- [38] S. G. Finlayson, J. D. Bowers, J. Ito, J. L. Zittrain, A. L. Beam, and I. S. Kohane, “Adversarial attacks on medical machine learning,” *Science*, vol. 363, no. 6433, pp. 1287–1289, 2019.
- [39] Y. Li, H. Zhang, C. Bermudez, Y. Chen, B. A. Landman, and Y. Vorobeychik, “Anatomical context protects deep learning from adversarial perturbations in medical imaging,” *Neurocomputing*, 2019.
- [40] M. Paschali, W. Simson, A. Guha Roy, M. Ferjad Naeem, R. Göbl, C. Wachinger, and N. Navab, “Data Augmentation with Manifold Exploring Geometric Transformations for Increased Performance and Robustness,” *arXiv e-prints*, p. arXiv:1901.04420, Jan 2019.

- [41] Y. Huang, T. Würfl, K. Breininger, L. Liu, G. Lauritsch, and A. Maier, “Some investigations on robustness of deep learning in limited angle tomography,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2018, pp. 145–153.
- [42] R. Paul, M. Schabath, R. Gillies, L. Hall, and D. Goldgof, “Mitigating adversarial attacks on medical image understanding systems,” in *2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2020, pp. 1517–1521.
- [43] X. Ma, Y. Niu, L. Gu, Y. Wang, Y. Zhao, J. Bailey, and F. Lu, “Understanding adversarial attacks on deep learning based medical image analysis systems,” *arXiv preprint arXiv:1907.10456*, 2019.
- [44] X. Li and D. Zhu, “Robust detection of adversarial attacks on medical images,” in *2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2020, pp. 1154–1158.
- [45] Z. Zhang, Q. Liu, and Y. Wang, “Road extraction by deep residual U-Net,” *IEEE Geoscience and Remote Sensing Letters*, vol. 15, no. 5, pp. 749–753, 2018.
- [46] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards real-time object detection with region proposal networks,” in *Advances in neural information processing systems*, 2015, pp. 91–99.
- [47] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [48] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville, “Improved Training of Wasserstein GANs,” *ArXiv e-prints*, p. arXiv:1704.00028, 2017.
- [49] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, “Spectral normalization for generative adversarial networks,” *ArXiv e-prints*, p. arXiv:1802.05957, Feb. 2018.