
AT-GAN: A Generative Attack Model for Adversarial Transferring on Generative Adversarial Net

Xiaosen Wang, Kun He*

School of Computer Science and Technology
Huazhong University of Science and Technology
Wuhan 430074, China
{xiaosen, brooklet60}@hust.edu.cn

John E. Hopcroft

Department of Computer Science, Cornell University
Ithaca 14853, USA
jeh@cornell.edu

Abstract

Recent studies have discovered the vulnerability of Deep Neural Networks (DNNs) to *adversarial examples* that are imperceptible to humans but can easily fool DNNs. Existing methods for crafting adversarial examples are mainly based on *search in the neighborhood of the input*, such as adding small-magnitude perturbations to the original images so that the generated adversarial examples differ from the benign examples by a small matrix norm. In this work, we propose a new generative model called AT-GAN to estimate the distribution of adversarial examples so that it can directly generate adversarial examples from any random noise using Generative Adversarial Nets (GANs). The key idea is to transfer a pre-trained GAN to generate adversarial examples for the target classifier. Once trained, AT-GAN can generate plentiful, high quality and diverse adversarial examples very quickly. We evaluate AT-GAN in both semi-whitebox and black-box settings under typical defense methods on two benchmark datasets, MNIST and Fashion-MNIST. Empirical comparisons with existing attack baselines demonstrate that AT-GAN can achieve a higher attack success rate with a better efficiency.

1 Introduction

Deep Neural Networks (DNNs) have exhibited great performance in computer vision tasks in recent years (Krizhevsky et al., 2012; He et al., 2016). However, DNNs are found vulnerable to adversarial examples (Szegedy et al., 2014). Due to the robustness and security implications, the ways of generating adversarial examples are called *attacks* and there are two types of adversarial attacks, targeted and untargeted. *Targeted attacks* aim to generate adversarial examples that are classified as specific target classes, while *untargeted attacks* aim to generate adversarial examples that are classified incorrectly. Various algorithms (Yuan et al., 2017) have been proposed for generating adversarial examples, such as the Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2015) and Carlini-Wagner attack (C&W) (Carlini and Wagner, 2017). Adversarial examples can be used for the training to improve the model’s robustness, which is a popular and efficient defense method called *adversarial training* (Goodfellow et al., 2015; Kurakin et al., 2017; Song et al., 2019).

*The first two authors contribute equally.

In order to generate adversarial examples, most attack algorithms (Goodfellow et al., 2015; Carlini and Wagner, 2017) add imperceptible perturbation to the input based on gradient descent which means their generated adversarial examples are restricted by the original images. Xiao et al. (2018) propose to train a generator G using a benign image x as the input to generate the perturbation $G(x)$ so that $x + G(x)$ can fool the target model. However, their result is still restricted by the original images. Song et al. (2018) assume that the images generated by GAN contain adversarial examples, so they propose a new method that searches a noise input near an arbitrary noise vector using gradient descent of the target classifier for AC-GAN such that this input leads to adversarial example. Song et al. (2018) called their output *unrestricted adversarial examples* by claiming that their output is not limited to a benign image. However, their output is still limited to the noise input as they use gradient descent to search for a good noise in the neighborhood of the original noise. Because their method can not always succeed for any noise, they need to switch to another random noise as the input in case the current search fails. Besides, their method has a slow generation speed by involving hundreds of iterations on the gradient descent in order to find a good noise in the neighborhood.

In this work, we propose a new generative attack model called AT-GAN (Adversarial Transferring on Generative Adversarial Nets) to estimate the distribution of adversarial examples so as to generate *unrestricted adversarial examples* from random noise. We first normally train a Generative Adversarial Network (GAN) (Goodfellow et al., 2014) and then transfer the generator to attack the target classifier. Note that compared to Song et al. (2018), our output is *truly* unrestricted to the input because it has learned the distribution of adversarial examples. Once our generator is transferred from generating normal images to adversarial images, it can directly generate unrestricted adversarial examples with any random noise, leading to a high diversity. Also, we do not depend on iterations of the gradient method, so our generation is very fast.

To evaluate the effectiveness of our attack strategy, we use AT-GAN on several models to generate adversarial examples from random noise and compare our model with several other attack methods in both semi-whitebox and black-box settings. Then we apply typical defense methods (Goodfellow et al., 2015; Madry et al., 2017; Tramèr et al., 2018) to defend against these generated adversarial examples. Empirical results show that the adversarial examples generated by AT-GAN yields a higher attack success rate. Our main contributions are as follows:

- We use GAN to estimate the distribution of adversarial examples so as to generate unrestricted adversarial examples from random noise. This is different from most existing attack methods that focus on how to add crafted perturbation to the original image.
- We train a conditional generative network to directly produce adversarial examples, which does not rely on the gradient of the input. The generation process is very fast. This is different from the main stream of attacks based on optimization.
- As compared with the very few works using GAN for attack, which are still based on the searching near the neighborhood of the input, our generated images have a higher diversity because our method does not need a suitable noise as the input.
- Extensive empirical study on typical defense methods against adversarial examples shows that the proposed AT-GAN achieves a higher attack success rate than existing adversarial attacks on both semi-whitebox and black-box settings.

2 Related Work

In this section, we provide an overview of GANs and typical attack methods to generate adversarial examples. We also introduce several attacks based on GANs that are most related to our work. The typical defense methods based on adversarial training are described in Appendix A.

2.1 Generative Adversarial Nets (GANs)

A generative adversarial net (GAN) (Goodfellow et al., 2014) consists of two neural networks trained in opposition to each other. The generator G is optimized to estimate the data distribution and the discriminator D aims to distinguish fake samples from G and real samples from the training data. The objective of D and G can be formalized as a min-max value function $V(G, D)$:

$$\min_G \max_D V(G, D) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \quad (1)$$

The Conditional Generative Adversarial Net (CGAN) (Mirza and Osindero, 2014) is the first conditional version of GAN, which combines conditions from the input of both generator and discriminator. Radford et al. (2016) propose a Deep Convolutional Generative Adversarial Net (DCGAN), which implements GAN with convolutional networks and stabilizes the model during the training. Auxiliary Classifier GAN (AC-GAN) is another variant that extends GAN with some conditions by an extra classifier (Odena et al., 2017). The objective function of AC-GAN is as follows:

$$\min_G \max_D \min_C V(G, D, C) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z, y)))] \\ + \mathbb{E}_{x \sim p_{data}(x)} [\log(1 - C(x, y))] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - C(G(z, y), y))] \quad (2)$$

To make the GAN more trainable in practice, Arjovsky et al. (2017) proposed Wasserstein GAN (WGAN) which uses Wasserstein distance so that the loss function has more desirable properties. Gulrajani et al. (2017) introduce WGAN_GP (WGAN with gradient penalty) that performs better than WGAN in practice. Its objective function is formulated as follows:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}} [D(x)] - \mathbb{E}_{z \sim p_z} [D(G(z))] - \lambda \mathbb{E}_{\hat{x} \sim p_{\hat{x}}} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2] \quad (3)$$

2.2 Gradient based Attacks

There are three types of attacks regarding how the attacks access the model. The *white-box attack* can fully access the target model, while the *black-box attack* (Papernot et al., 2017) has no knowledge of the target model. Existing black-box attacks mainly focus on *transferability* (Liu et al., 2017; Bhagoji et al., 2017), in which an adversarial instance generated on one model could be transferred to attack another model. A third type of *semi-whitebox attack* was recently proposed by Xiao et al. (2018). Semi-whitebox attack needs the logits output during training but afterwards can generate adversarial examples without accessing the target model.

We will introduce three typical adversarial attack methods. Here the components of all adversarial examples are clipped in $[0, 1]$.

Fast Gradient Sign Method (FGSM). FGSM (Goodfellow et al., 2015) adds perturbation in the direction of the gradient of the training loss J on the input x to generate adversarial examples.

$$x^{adv} = x + \epsilon \cdot \text{sign}(\nabla_x J(\theta, x, y)) \quad (4)$$

Here y is the true label of a sample x , θ is the model parameter and ϵ specifies the ℓ_∞ distortion between x and x^{adv} .

Projected Gradient Descent (PGD). The PGD adversary (Madry et al., 2017) is a multi-step variant of FGSM, which applies FGSM iteratively for k times with a budget α .

$$x^{adv_{t+1}} = \text{clip}(x^{adv_t} + \alpha \text{sign}(\nabla_x J(\theta, x^{adv_t}, y)), x^{adv_t} - \epsilon, x^{adv_t} + \epsilon) \quad \text{where } x^{adv_0} = x \\ x^{adv} = x^{adv_k} \quad (5)$$

Here $\text{clip}(x, p, q)$ forces its input x to reside in the range of $[p, q]$.

Rand FGSM (R+FGSM). R+FGSM (Tramèr et al., 2018) first applies a small random perturbation on the benign image with parameter α (where $\alpha < \epsilon$), then it uses FGSM to generate the adversarial example based on the perturbed image.

$$x^{adv} = x' + (\epsilon - \alpha) \cdot \text{sign}(\nabla_{x'} J(\theta, x', y)) \quad \text{where } x' = x + \alpha \cdot \text{sign}(\mathcal{N}(\mathbf{0}^d, \mathbf{I}^d)) \quad (6)$$

2.3 GAN based Attacks

In this work we propose to use GAN to generate adversarial examples. In the literature, there are only a few attack methods based on GANs, such as AdvGAN (Xiao et al., 2018) and unrestricted adversarial attacks (Song et al., 2018).

AdvGAN. Xiao et al. (2018) propose to train an AdvGAN to take a benign image x as the input and generate a perturbation $G(x)$ which aims to make the target model f classify $x + G(x)$ in target class y_t . The objective function is formulated as:

$$\min_G \max_D V(G, D, f) = \mathbb{E}_{x \sim p_x(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(x + G(x)))] \\ + \mathbb{E}_{x \sim p_x(x)} [1 - f(x + G(x), y_t)] + \mathbb{E}_{x \sim p_x(x)} [\max(0, \|G(x)\|_2 - c)] \quad (7)$$

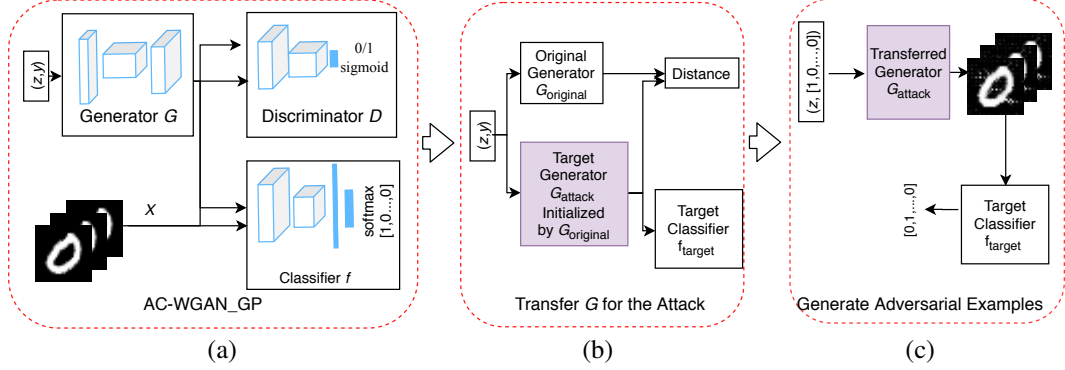


Figure 1: The architecture of AT-GAN. The first training stage of AT-GAN is similar to that of AC-GAN. After G is trained, we regard G as the original model $G_{original}$ and copy $G_{original}$ as the initial attack model G_{attack} . We then transfer G_{attack} according to the target classifier to be attacked. After the second stage of training, AT-GAN can generate adversarial examples by G_{attack} .

Unrestricted adversarial attacks. Song et al. (2018) propose to search a noise input z^* near an arbitrary noise vector z^0 for AC-GAN so as to produce an adversarial example $G(z^*, y_s)$ for the target model f with an extra classifier C . The objective function can be written as:

$$z^* = \arg \min_z \left\{ \frac{\lambda_1}{m} \sum_{i=1}^m \max(|z_i - z_i^0| - \epsilon, 0) - \lambda_2 \log C(G(z^*, y_s), y_s) - \log f(G(z^*, y_s), y_t) \right\} \quad (8)$$

Our proposed AT-GAN is implemented based on the combination of AC-GAN and WGAN_GP (AC-WGAN_GP), but we have a very different objective for the training. We aim to learn the distribution of adversarial examples so that we could generate plentiful and diverse images that are visually realistic but misclassified by the target classifier. See details in the next section.

3 The proposed AT-GAN

In order to generate adversarial examples from random noise, we seek a generator G_{attack} to fit the distribution of adversarial examples.

$$f_{target}(G_{attack}(z, y)) = y_t \quad s. t. \quad G_{attack}(z, y) \in S_y, y \neq y_t \quad (9)$$

Here f_{target} is the target classifier to be attacked, z is a random noise, y is the true label of the generated image, y_t is the target label, and S_y represents the set of images with label y .

The above objective function is hard to solve, so we propose a new model called AT-GAN (Adversarial Transferring on Generative Adversarial Net). The architecture of AT-GAN is illustrated in Figure 1. There are two training stages. We first train the GAN model to get a generator $G_{original}$ (See details in Appendix B), then we transfer $G_{original}$ to attack f_{target} .

3.1 Transfer the Generator for Attack

After the original generator G is trained, we transfer the generator to learn the distribution of adversarial examples in order to attack the target network. As illustrated in Figure 1 (b), there are three neural networks, the original generator $G_{original}$, the attack generator G_{attack} to be transferred that has the same set of weights as $G_{original}$ in the beginning and a classifier f_{target} to be attacked. Then Eq. 9 could be rewritten as follows:

$$G_{attack} = \arg \min_G \|G_{original}(z, y) - G(z, y)\|_p \quad s. t. \quad g(f_{target}(G(z, y))) = y \quad (10)$$

where $g(x)$ is the inverse function of $h(y) = y_t$, $\|\cdot\|_p$ is the ℓ_p norm. We use ℓ_2 norm in experiments.

Based on Eq. 10, we construct the loss function by two components. L_a aims to assure that f_{target} yields the target label y_t .

$$L_a(z, y) = H(g(f_{target}(G_{attack}(z, y))), y_t) \quad (11)$$

And L_d aims to assure that G_{attack} generate realistic examples.

$$L_d(z, y) = \|G_{original}(z, y) + P - G_{attack}(z, y)\|_p \quad (12)$$

where P is a Gaussian noise constrained by both l_0 and l_∞ norm. In practice, there is a small difference from Eq. 10 as we add a Gaussian noise to smooth the generated adversarial example.

Thus, the total loss for transferring G_{attack} can be written as:

$$L(z, y) = \alpha L_a(z, y) + \beta L_d(z, y) \quad (13)$$

Here α and β are hyperparameters to control the training process. For the untargeted attack, we just replace y_t in L_a with $\max_{y \neq y_t} C(G_{attack}(z, y))$ where $C(\cdot)$ is the logits of the target classifier. Note that when $\alpha = 1$ and $\beta \rightarrow \infty$, the objective function will be similar to FGSM (Goodfellow et al., 2015) or R+FGSM (Tramèr et al., 2018).

3.2 Learn the Distribution of Adversarial Examples

Given a distribution $p(x)$ in space \mathcal{X} , we construct another distribution $q(x)$ by choosing a point $p'(x)$ in the ϵ -neighborhood of $p(x)$ for any $x \in \mathcal{X}$. Obviously, when $p'(x)$ is close enough to $p(x)$, $q(x)$ has almost the same distribution as $p(x)$.

Lemma 1. *Given two distributions P and Q with probability density function $p(x)$ and $q(x)$ in a space \mathcal{X} , if there is a ϵ that satisfies $\|q(x) - p(x)\| < \epsilon$ for any $x \in \mathcal{X}$, we could get $KL(P\|Q) \rightarrow 0$ when $\epsilon \rightarrow 0$.*

The proof is in Appendix C.1. Eq. 12 aims to constrain the image generated by G_{attack} in the ϵ -neighborhood of $G_{original}$. Under the ideal condition that Eq. 12 guarantees $G_{attack}(z)$ is close enough to $G_{original}(z)$ in high dimension for any input noise z , the distribution of AT-GAN almost coincides that of WGAN_GP, aka $p_a \approx p_g$.

Samangouei et al. (2018) prove that the global optimum of WGAN is $p_g = p_{data}$, when the set $\{x | p_g(x) \neq p_{data}(x)\}$ has zero Lebesgue-measure. Similarly, we could show that the optimum of WGAN_GP is $p_g = p_{data}$ under the same condition. The proof is in Appendix C.2.

Therefore, under ideal conditions, we conclude $p_a \approx p_g = p_{data}$, which means p_a has almost the same distribution as the real data. However, the images generated by G_{attack} are adversarial examples. Thus, the distribution p_a learned by G_{attack} is the distribution of adversarial examples.

4 Experiments

We evaluate the proposed AT-GAN² with several neural models on two benchmark datasets. The results demonstrate that for both semi-whitebox and black-box attacks, AT-GAN could achieve a higher attack success rate than the state-of-art baselines and AT-GAN indeed learns a good distribution of adversarial examples close to the real data distribution. Besides, AT-GAN is very fast.

4.1 Experimental Setup

All experiments are done on a single Titan X GPU. We test four neural network models, Models A, B, C, and D. Their details on neural network architecture and hyperparameters are in Appendix D.

Baselines. We compare AT-GAN with typical state-of-art attacks, namely FGSM, PGD, R+FGSM and Song et al. (2018). And we evaluate their attack performance on several defense methods, namely adversarial training Goodfellow et al. (2015), ensemble adversarial training Tramèr et al. (2018) and iterative adversarial training Madry et al. (2017).

Datasets. MNIST (LeCun and Cortes, 1998) is a dataset of hand written digit number from 0 to 9. Fashion-MNIST (Xiao et al., 2017) is similar to MNIST with ten classes of fashion clothes. The corresponding class labels from 0 to 9 are t-shirt/top, trouser, pullover, dress, coat, sandal, shirt, sneaker, bag and ankle boot, respectively. For both datasets, we normalize the pixel value range converting $[0, 255]$ to $[0, 1]$.

Attack success rate. For the attack performance, we use attack success rate, defined as the fraction of samples that meets the goal of the adversary: $f(x_{adv}) \neq y$ for untargeted attacks and $f(x_{adv}) = y_t$ for targeted attacks with a target y_t .

²The codes will be public after review.

4.2 Comparison on Attacks

We train the four neural Models A, B, C, and D by normal training, and the above three adversarial training methods, respectively. The classification accuracy of each model on the original test set is shown in Table 1. Models A, B and C achieve an accuracy rate greater than 99% on MNIST and 89.5% on Fashion-MNIST. Model D gains an accuracy rate between 97.3% and 98.5% on MNIST while between 85.2% and 89.4% on Fashion-MNIST, which is slightly lower as it does not contain any convolutional layer. We compare the attack performance of AT-GAN with other attack methods on the generation efficiency and classification accuracy on adversarial instances with these models.

Table 1: The classification accuracy on test set for various models. Nor.: Normal training, Adv.: Adversarial training, Ens.: Ensemble adversarial training, Iter. Adv.: Iterative adversarial training.

	MNIST(%)				Fashion-MNIST(%)			
	A	B	C	D	A	B	C	D
Nor.	99.1	99.0	99.0	97.3	91.8	90.2	92.2	85.2
Adv.	99.1	99.1	99.0	97.4	91.7	90.7	91.9	86.3
Ens.	99.2	99.0	99.0	97.4	92.0	89.6	92.1	85.5
Iter. Adv.	99.1	99.1	99.0	98.5	89.5	90.7	90.0	89.4

4.2.1 Comparison on Attack Efficiency

We first test the efficiency of each attack method using Model A on MNIST data. The average time of generating 1000 adversarial examples is listed in Table 2. Among the five attack methods, AT-GAN is the fastest as it could generate adversarial examples without using the target classifier.

Table 2: Comparison on the example generating time, measured by generating 1000 adversarial instances using Model A on MNIST.

	FGSM	PGD	R+FGSM	Song's	AT-GAN
Generating time	0.3s	1.8s	0.4s	>15min	0.2s

4.2.2 Comparison on the Generated Adversarial Examples

The adversarial examples generated by different methods for Model A are shown in Figure 2 and 3. On MNIST dataset, AT-GAN and Song's method (Song et al., 2018) generate much higher quality examples than other attacks, and AT-GAN generates slightly more realistic images than Song's method, for example on "0" and "3". On Fashion-MNIST dataset, some adversarial examples generated by Song's method are not realistic to human eyes, for example on "t-shirt/top (0)" and "sandal (5)". AT-GAN can generate very realistic images as Eq. 12 forces the adversarial examples to be close enough to the images generated by the original generator so they are realistic. Other attacks tend to perturb the images and some images are not clear enough for human eyes.



Figure 2: Adversarial examples generated by different methods using Model A on MNIST.

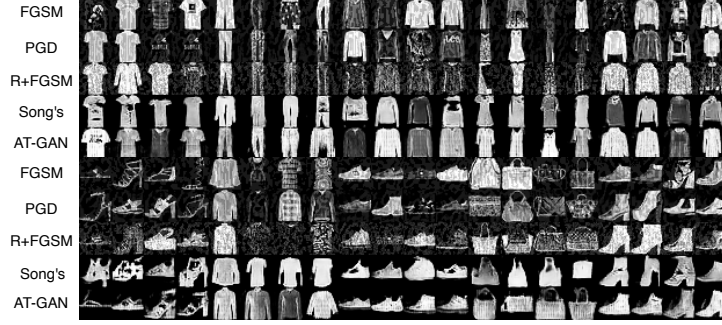


Figure 3: Adversarial examples generated by different methods using Model A on Fashion-MNIST.

4.2.3 Comparison on Semi-whitebox or Whitebox Attacks

As AT-GAN does not need to know the target model after it is trained, we generate adversarial instances by AT-GAN under semi-whitebox setting. For comparison, we also generate adversarial instances by the four baseline attack methods under white-box setting. We test on the four neural models by normal training and several adversarial trainings as the defense methods respectively. The attack success rates are listed in Table 3.

On MNIST dataset, AT-GAN is the best on all defenses. Only for Model A, B and C by normal training, PGD achieves the highest attack success rate of 100%, and AT-GAN gains the second highest attack success rate of over 99%. For all other cases, AT-GAN achieves the highest attack success rate.

On Fashion-MNIST dataset, AT-GAN is the best on average. PGD achieves the highest attack success rate on Model A with normal training and Model C with normal training and ensemble adversarial training. Song et al. (2018) achieves the highest attack success rate on Model A with various adversarial trainings. On these cases, AT-GAN almost achieves the second highest attack success rate close to the highest one. For all other cases, AT-GAN achieves the highest attack success rate.

Moreover, as compared with the normal training, AT-GAN has a good attack success rate on these defense methods while the attack success rate of the baselines clearly decay. This may be due to AT-GAN trying to estimate the distribution of adversarial examples so the adversarial training defenses do not have much impact on AT-GAN.

Table 3: Attack success rate of adversarial examples generated by AT-GAN in semi-whitebox attack and the baselines in white-box attack under several defenses on MNIST. On each model, the highest or second highest attack success rates are highlighted in **bold** or underline.

Model	Defense	MNIST(%)					Fashion-MNIST(%)				
		FGSM	PGD	R+FGSM	Song's	AT-GAN	FGSM	PGD	R+FGSM	Song's	AT-GAN
A	Nor.	68.8	100.0	77.9	57.3	<u>98.7</u>	82.7	99.9	95.5	89.2	<u>96.1</u>
	Adv.	2.9	<u>92.6</u>	28.4	34.3	97.5	14.6	82.6	63.3	93.1	<u>92.7</u>
	Ens.	8.7	<u>85.1</u>	20.3	36.6	96.7	36.0	90.8	68.4	95.9	<u>95.4</u>
	Iter. Adv.	4.0	5.9	2.6	<u>40.4</u>	91.4	23.2	30.2	37.8	95.9	<u>93.5</u>
B	Nor.	88.0	100.0	96.5	39.2	<u>99.5</u>	82.3	<u>96.0</u>	90.9	80.1	<u>98.5</u>
	Adv.	9.5	<u>42.4</u>	7.1	22.9	97.7	24.2	69.9	74.2	<u>81.6</u>	91.1
	Ens.	18.0	<u>98.2</u>	42.1	34.5	99.3	30.0	<u>95.3</u>	81.3	<u>81.8</u>	93.6
	Iter. Adv.	9.2	<u>36.2</u>	4.6	31.8	95.6	23.8	<u>34.7</u>	28.2	<u>72.0</u>	91.6
C	Nor.	70.7	100	81.6	42.4	<u>99.3</u>	82.7	100.0	98.4	94.3	98.0
	Adv.	4.3	<u>76.9</u>	7.2	37.5	95.8	11.0	88.3	76.5	88.8	88.9
	Ens.	7.8	<u>96.4</u>	19.9	41.7	96.9	47.8	99.9	71.3	88.6	<u>93.9</u>
	Iter. Adv.	4.7	9.6	2.9	<u>31.5</u>	90.0	22.3	28.8	31.2	<u>72.7</u>	91.6
D	Nor.	89.6	91.7	<u>93.8</u>	46.6	99.9	77.2	85.4	<u>88.3</u>	70.3	99.9
	Adv.	23.5	<u>96.8</u>	76.2	21.7	99.9	33.8	54.7	45.3	<u>55.1</u>	99.4
	Ens.	34.6	99.5	51.6	34.2	99.5	47.5	<u>76.8</u>	62.1	65.7	99.1
	Iter. Adv.	49.8	<u>81.1</u>	25.3	18.1	99.7	22.9	30.2	33.6	<u>75.3</u>	93.1

4.2.4 Transferability for Black-box Attacks

For transferability evaluation, we use **Model A** to perform black-box attack against the target **Model C**. The attack success rate is shown in Table 4. On MNIST dataset, adversarial examples generated by PGD transfer best for normal training, while adversarial instances generated by AT-GAN achieve much higher attack success rate for adversarial training as the defense method. On Fashion-MNIST dataset, PGD achieves the highest attack rate on models with normal training and ensemble adversarial training, and Song et al. (2018) is the best on the other two adversarial trainings. However, this does not mean that AT-GAN has a lower transferability because the adversarial examples generated by other attacks on Fashion-MNIST are not realistic as that of AT-GAN, and less realistic images lead to a higher attack rate.

Table 4: Attack success rate of adversarial examples generated by different black-box adversarial strategies on **Model C** based on examples generated on **Model A**.

	MNIST(%)				Fashion-MNIST(%)			
	Nor.	Adv.	Ens.	Iter. Avd.	Nor.	Adv.	Ens.	Iter. Avd.
FGSM	46.7	4.2	1.7	4.6	68.9	23.1	20.8	14.8
PGD	97.5	6.5	4.1	4.1	84.7	27.6	39.6	14.6
R+FGSM	82.3	6.7	4.8	4.1	21.2	32.1	17.5	26.3
Song's	5.7	3.8	3.5	3.2	38.1	32.9	31.3	31.4
AT-GAN	65.3	24.6	27.9	17.2	58.0	22.7	32.0	15.23

4.3 Exploration on the Distribution of Adversarial Examples

To identify that AT-GAN can learn a distribution of adversarial examples which is close to the distribution of real data, we use t-SNE (Maaten and Hinton, 2008) on 5,000 real images sampled from test set and 5,000 generated adversarial examples to illustrate the distribution in 2 dimensions. If the adversarial examples have a different distribution from that of the real data, then t-SNE cannot deal well with them and the result will be in chaos. The results are illustrated in Figure 4.3. On both datasets of MNIST and Fashion-MNIST, the result of AT-GAN is closest to that of the test set. It indicates that AT-GAN indeed learns the distribution that is closest to the distribution of the real data.

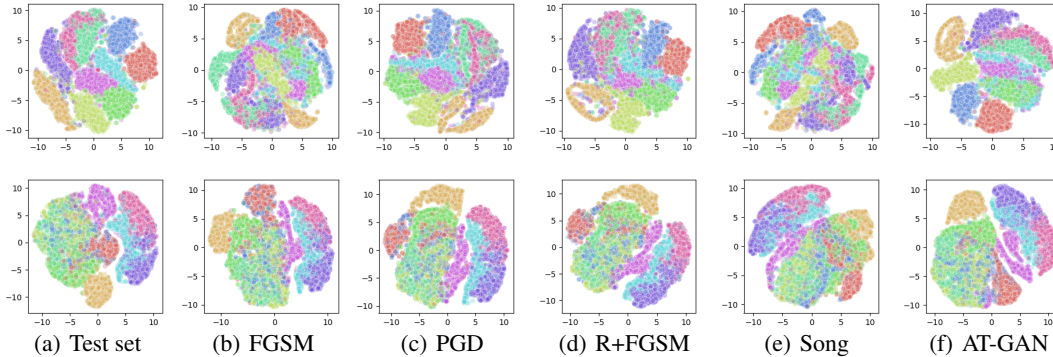


Figure 4: The t-SNE result for test set and adversarial examples generated by different methods on MNIST (top) and Fashion-MNIST (bottom). For (a), we use 10,000 sampled real images in test set. For (b) to (f), we use 5,000 sampled images in test set and 5,000 adversarial examples by different attacks. The position of each class is random due to the property of t-SNE.

5 Conclusion

We propose a new attack model called AT-GAN that transfers generative adversarial nets (GANs) for adversarial attacks. Different from existing attacks that add small perturbations to the input images, AT-GAN tries to explore the distribution of adversarial instances so as to directly generate the adversarial examples from any random noise. In this way, the generated adversarial examples are not limited to any natural images. Also, compared with the pioneer work using GAN (Song

et al., 2018) that do iterations of gradient descent to search for a good noise in the neighborhood of an original random noise such that the corresponding output of GAN is an adversarial example, our model is a generative model that tries to learn the distribution of the adversarial examples. Once AT-GAN is trained, it can generate adversarial images and the output is not limited to the input noise. Experiments show that AT-GAN is very fast and can generate plenty of adversarial instances that look more realistic to human eyes, AT-GAN yields a higher attack success rate under various adversarial training defenses for semi-whitebox as well as black-box attack settings, and AT-GAN can learn a distribution of adversarial examples that is very close to the distribution of the real data.

Acknowledgments

We thank Chuan Guo for helpful discussions and suggestions on our work.

References

- Arjovsky, M., Chintala, S., and Bottou, L. (2017). Wasserstein GAN. In *arXiv preprint arXiv:1701.07875*.
- Athalye, A., Carlini, N., and Wagner, D. (2018). Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples. In *International Conference on Machine Learning*.
- Bhagoji, A. N., He, W., Li, B., and Song, D. (2017). Exploring the Space of Black-box Attacks on Deep Neural Networks. In *arXiv preprint arXiv:1703.09387*.
- Buckman, J., Roy, A., Raffel, C., and Goodfellow, I. (2018). Thermometer Encoding: One Hot Way To Resist Adversarial Examples. In *International Conference on Learning Representations*.
- Carlini, N. and Wagner, D. (2017). Towards Evaluating the Robustness of Neural Networks. In *IEEE Symposium on Security and Privacy*, pages 39–57.
- Chen, X., Duan, Y., Houthoofd, R., Schulman, J., Sutskever, I., and Abbeel, P. (2016). InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets. In *Neural Information Processing Systems*.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative Adversarial Nets. In *International Conference on Learning Representations*.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. (2015). Explaining and Harnessing Adversarial Examples. In *International Conference on Learning Representations*.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A. (2017). Improved Training of Wasserstein GANs. In *Neural Information Processing Systems*.
- Guo, C., Rana, M., Cisse, M., and van der Maaten, L. (2018). Countering adversarial images using input transformations. In *International Conference on Learning Representations*.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep Residual Learning for Image Recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). ImageNet Classification with Deep Convolutional Neural Networks. In *Neural Information Processing Systems*.
- Kurakin, A., Goodfellow, I., and Bengio, S. (2017). Adversarial Machine Learning at Scale. In *International Conference on Learning Representations*.
- LeCun, Y. and Cortes, C. (1998). The MNIST database of handwritten digits.
- Liao, F., Liang, M., Dong, Y., Pang, T., Hu, X., and Zhu, J. (2018). Defense against Adversarial Attacks Using High-Level Representation Guided Denoiser. In *IEEE Conference on Computer Vision and Pattern Recognition*.

- Liu, Y., Chen, X., Liu, C., and Song, D. (2017). Delving into Transferable Adversarial Examples and Black-box Attacks. In *International Conference on Learning Representations*.
- Maaten, L. and Hinton, G. (2008). Visualizing Data using t-SNE. In *Neural Information Processing Systems*.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. (2017). Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference on Learning Representations*.
- Metzen, J. H., Genewein, T., Fischer, V., and Bischoff, B. (2017). On Detecting Adversarial Perturbations. In *International Conference on Learning Representations*.
- Mirza, M. and Osindero, S. (2014). Conditional Generative Adversarial Nets. In *arXiv preprint arXiv:1411.1784*.
- Odena, A., Olah, C., and Shlens, J. (2017). Conditional Image Synthesis With Auxiliary Classifier GANs. In *International Conference on Machine Learning*.
- Papernot, N., McDaniel, P., Goodfellow, I., Jha, S., Celik, Z. B., and Swami, A. (2017). Practical Black-Box Attacks against Machine Learning. In *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security*, page 506–519.
- Radford, A., Metz, L., and Chintala, S. (2016). Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. In *International Conference on Learning Representations*.
- Samangouei, P., Kabkab, M., and Chellappa, R. (2018). Defense-GAN: Protecting Classifiers Against Adversarial Attacks Using Generative Models. In *International Conference on Learning Representations*.
- Shen, S., Jin, G., Gao, K., and Zhang, Y. (2017). APE-GAN: Adversarial Perturbation Elimination with GAN. In *arXiv preprint arXiv:1707.05474*.
- Song, C., He, K., Wang, L., and Hopcroft, J. E. (2019). Improving the Generalization of Adversarial Training with Domain Adaptation. In *International Conference on Learning Representations*.
- Song, Y., Shu, R., Kushman, N., and Ermon, S. (2018). Constructing Unrestricted Adversarial Examples with Generative Models. In *Neural Information Processing Systems*.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. (2014). Intriguing properties of neural networks. In *Neural Information Processing Systems*.
- Tramèr, F., Kurakin, A., Papernot, N., Goodfellow, I., Boneh, D., and McDaniel, P. (2018). Ensemble Adversarial Training: Attacks and Defenses. In *International Conference on Learning Representations*.
- Xiao, C., Li, B., Zhu, J.-Y., He, W., Liu, M., and Song, D. (2018). Generating Adversarial Examples with Adversarial Networks. In *International Joint Conferences on Artificial Intelligence*.
- Xiao, H., Rasul, K., and Vollgraf, R. (2017). Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. In *arXiv preprint arXiv:1708.07747*.
- Yuan, X., He, P., Zhu, Q., and Li, X. (2017). Adversarial Examples: Attacks and Defenses for Deep Learning. In *arXiv preprint arXiv:1712.07107*.

Appendix

A Adversarial training based Defenses

There are many defense strategies, such as detecting adversarial perturbation (Metzen et al., 2017), obfuscating gradients (Buckman et al., 2018; Guo et al., 2018) and eliminating perturbation (Shen et al., 2017; Liao et al., 2018), among which adversarial training is the most effective method (Athalye et al., 2018). We list several adversarial training methods as follows.

Adversarial training. Goodfellow et al. (2015) first introduce the method of adversarial training, where the standard loss function f for a neural network is modified as:

$$\tilde{J}(\theta, x, y) = \alpha J_f(\theta, x, y) + (1 - \alpha) J_f(\theta, x^{adv}, y) \quad (14)$$

Here y is the true label of a sample x and θ is the model's parameter. The modified objective is to make the neural network more robust by penalizing it to count for adversarial samples. During the training, the adversarial samples are computed with respect to the current status of the network. Taking FGSM for example, the loss function could be written as:

$$\tilde{J}(\theta, x, y) = \alpha J_f(\theta, x, y) + (1 - \alpha) J_f(\theta, x + \epsilon \text{sign}(\nabla_x J(\theta, x, y)), y) \quad (15)$$

Ensemble adversarial training. Tramèr et al. (2018) propose an ensemble adversarial training method, in which DNN is trained with adversarial examples transferred from a number of fixed pre-trained models.

Iterative adversarial training. Madry et al. (2017) propose to train a DNN with adversarial examples generated by iterative methods such as PGD.

B Train WGAN_GP to Get the Original Generator

Figure 1 (a) illustrates the overall architecture of the original AC-WGAN_GP. There are three neural networks: a generator G , a discriminator D and a classifier f . The generator G takes a random noise z and a label y as the inputs and generates an image $G(z, y)$. It aims to generate an image $G(z, y)$ that is indistinguishable to discriminator D and makes the classifier f to output a label y . The loss function of G can be formulated as:

$$L_G(z, y) = -\mathbb{E}_{z \sim p_z(z)}[D(G(z))] + \mathbb{E}_{z \sim p_z(z)} H(f(G(z, y)), y) \quad (16)$$

Here $H(a, b)$ is the entropy between a and b . The discriminator D takes the training data x or the generated data $G(z, y)$ as input and tries to distinguish them. The loss function of D with a penalty on the gradient norm for random samples $\hat{x} \sim p_{\hat{x}}$ can be formulated as:

$$L_D(x, z, y) = \mathbb{E}_{x \sim p_{data}(x)}[D(x)] - \mathbb{E}_{z \sim p_z(z)}[D(G(z, y))] + \lambda \mathbb{E}_{\hat{x} \sim p_{\hat{x}}(\hat{x})}[(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2] \quad (17)$$

The classifier f takes the train data x or the generated data $G(z, y)$ as the input and predicts the corresponding label. There is no difference from other classifiers and the model is trained only on the training data. The loss of the classifier is:

$$L_f(x, y) = \mathbb{E}_{x \sim p_{data}(x)} H(f(x), y) \quad (18)$$

The goal of this stage is to train a generator G that could output realistic samples and estimate the data distribution properly so that we could later on transfer the generator to generate adversarial examples. So you could train a generator G using other GANs as long as the generator could learn a good distribution of the real data.

C Proofs that AT-GAN can Learn the Distribution of Adversarial Examples

C.1 AT-GAN has almost the same distribution as that of the original generator.

For two distributions P and Q with probability density function $p(x)$ and $q(x)$, we could get $q(x) = p(x) + f(x)$ where $\|f(x)\| < \epsilon$.

$$\begin{aligned}
KL(P\|Q) &= \int p(x) \log \frac{p(x)}{q(x)} dx \\
&= \int p(x) \log p(x) dx - \int p(x) \log q(x) dx \\
&= \int (q(x) - f(x)) \log p(x) dx - \int (q(x) - f(x)) \log q(x) dx \\
&= \int q(x) \log p(x) dx - \int q(x) \log q(x) dx - \int f(x) \log p(x) dx + \int f(x) \log q(x) dx \\
&= \int f(x) \log \frac{q(x)}{p(x)} dx - KL(Q\|P) \\
&\leq \int \epsilon \log(1 + \frac{\epsilon}{p(x)}) dx
\end{aligned} \tag{19}$$

Obviously, when $\epsilon \rightarrow 0$, we could get $\int \epsilon \log(1 + \frac{\epsilon}{p(x)}) dx \rightarrow 0$, which means $KL(P\|Q) \rightarrow 0$.

C.2 Global Optimality of $p_g = p_{data}$ for WGAN_GP

We refers to Samangouei et al. (2018) to prove this property. The WGAN_GP min-max loss is given by:

$$\begin{aligned}
\min_G \max_D V(D, G) &= \mathbb{E}_{x \sim p_{data}}[D(x)] - \mathbb{E}_{z \sim p_z(z)}[D(G(z))] - \lambda \mathbb{E}_{\hat{x} \sim p_{\hat{x}}(\hat{x})}[(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2] \\
&= \int_x p_{data}(x) D(x) dx - \int_z p_z(z) D(G(z)) dz - \lambda \int_{\hat{x}} p_{\hat{x}}(\hat{x}) [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2] d\hat{x} \\
&= \int_x [p_{data}(x) - p_g(x)] D(x) dx - \lambda \int_{\hat{x}} p_{\hat{x}}(\hat{x}) [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2] d\hat{x}
\end{aligned} \tag{20}$$

For a fixed G , the optimal discriminator D which maximizes $V(D, G)$ should be:

$$D_G^*(x) = \begin{cases} 1 & \text{if } p_{data}(x) \geq p_g(x) \\ 0 & \text{otherwise} \end{cases} \tag{21}$$

According to Eq. 20 and Eq. 21, we could get:

$$\begin{aligned}
V(D, G) &= \int_x [p_{data}(x) - p_g(x)] D(x) dx - \lambda \int_{\hat{x}} p_{\hat{x}}(\hat{x}) [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2] d\hat{x} \\
&= \int_{\{x | p_{data}(x) \geq p_g(x)\}} (p_{data}(x) - p_g(x)) dx - \lambda \int_{\hat{x}} p_{\hat{x}}(\hat{x}) d\hat{x} \\
&= \int_{\{x | p_{data}(x) \geq p_g(x)\}} (p_{data}(x) - p_g(x)) dx - \lambda
\end{aligned} \tag{22}$$

Let $\mathcal{X} = \{x | p_{data}(x) \geq p_g(x)\}$, to minimize Eq. 22, we need to set $p_{data}(x) = p_g(x)$ for any $x \in \mathcal{X}$. Then, since both p_g and p_{data} integrate to 1, we could get:

$$\int_{\mathcal{X}^c} p_g(x) dx = \int_{\mathcal{X}^c} p_{data}(x) dx \tag{23}$$

However, this contradicts Eq. 21 for $p_{data}(x) < p_g(x)$ and $x \in \mathcal{X}^c$, unless $\mu(\mathcal{X}^c) = 0$ where μ is the Lebesgue measure. This concludes the proof.

D More Experimental Details

We describe the details about the experiments, including attack hyperparamters and model architectures, and we provide more property of the generated adversarial examples.

D.1 Hyperparameters

The hyperparameters used in the experiments are described in Table 5.

Table 5: The hyperparamters.

Attack	Hyperparameters for attacks		Norm
	MNIST	Fashion-MNIST	
FGSM	$\epsilon = 0.3$	$\epsilon = 0.1$	ℓ_∞
PGD	$\epsilon = 0.3, \alpha = 0.075, \text{epochs} = 20$	$\epsilon = 0.1, \alpha = 0.01, \text{epochs} = 20$	ℓ_∞
R+FGSM	$\epsilon = 0.3, \alpha = 0.15$	$\epsilon = 0.2, \alpha = 0.1$	ℓ_∞
Song’s	$\lambda_1 = 100, \lambda_2 = 0, \text{epochs} = 1000$	$\lambda_1 = 100, \lambda_2 = 0, \text{epochs} = 1000$	N/A
AT-GAN	$\alpha = 2, \beta = 1, \text{epochs} = 500$	$\alpha = 2, \beta = 1, \text{epochs} = 1000$	N/A

D.2 The Architecture of Models

We describe our neural network architectures used in experiments. The abbreviations for components in the network are described in Table 6. The architecture of WGAN_GP with auxiliary classifier is shown in Table 7. The details of model A through D are described in Table 8 and the generator and discriminator are the same as in Chen et al. (2016). Both **Model A** and **Model C** have convolutional layers and fully connected layers. The difference is only on the size and number of convolutional filters. **Model B** uses dropout as its first layer and adopts a bigger covolutional filter so it has less number of parameters. **Model D** is a fully connected neural network with the least number of parameters and its accuracy will be lower than others because there is no convolutional layers.

Table 6: The meaning of abbreviations.

Abbreviation	Meaning
$\text{Conv}(m, k \times k)$	A convolutional layer with m filters, with filter size k
$\text{DeConv}(m, k \times k)$	A transposed convolutional layer with m filters, with filters size k
$\text{Dropout}(\alpha)$	A dropout layer with probability α
$\text{FC}(m)$	A fully connected layer with m outputs
Sigmoid	The sigmoid activation function
Relu	The Rectified Linear Unit activation function
$\text{LeakyRelu}(\alpha)$	The Leaky version of a Rectified Linear Unit with parameter α

Table 7: The architecture of WGAN_GP with auxiliary classifier.

Generator	Discriminator	Classifier
FC(1024) + Relu	Conv(64, 4×4) + LeakyRelu(0.2)	Conv(32, 3×3) + Relu
FC($7 \times 7 \times 128$) + Relu	Conv(128, 4×4) + LeakyRelu(0.2)	pooling(2, 2)
DeConv(64, 4×4) + Sigmoid	FC(1024) + LeakyRelu(0.2)	Conv(64, 3×3) + Relu
DeConv(1, 4×4) + Sigmoid	FC(1) + Sigmoid	pooling(2, 2)
		FC(1024)
		Dropout(0.4)
		FC(10) + Softmax

Table 8: The architectures of the **Models A** through **D** we used for classification. The number in parentheses in the title is the number of parameters for each model.

Model A (3,382,346)	Model B (710,218)	Model C (4,795,082)	Model D (509,410)
Conv(64, 5×5)+Relu	Dropout(0.2)	Conv(128, 3×3)+Relu	$\begin{bmatrix} \text{FC}(300)+\text{Relu} \\ \text{Dropout}(0.5) \end{bmatrix} \times 4$ FC(10) + Softmax
Conv(64, 5×5)+Relu	Conv(64, 8×8)+Relu	Conv(64, 3×3)+Relu	
Dropout(0.25)	Conv(128, 6×6)+Relu	Dropout(0.25)	
FC(128)+Relu	Conv(128, 5×5)+Relu	FC(128)+Relu	
Dropout(0.5)	Dropout(0.5)	Droopout(0.5)	
FC(10)+Softmax	FC(10)+Softmax	FC(10)+Softmax	

E Target Attack Examples of AT-GAN

We show some adversarial examples generated by AT-GAN with target attack. The results are illustrated in Fig. 5 and Fig. 6. Instead of adding perturbation to the original images, AT-GAN transfers the generator so that the adversarial instance would not have totally the same shape of the initial examples generated by the original GAN, as shown in the diagonal.

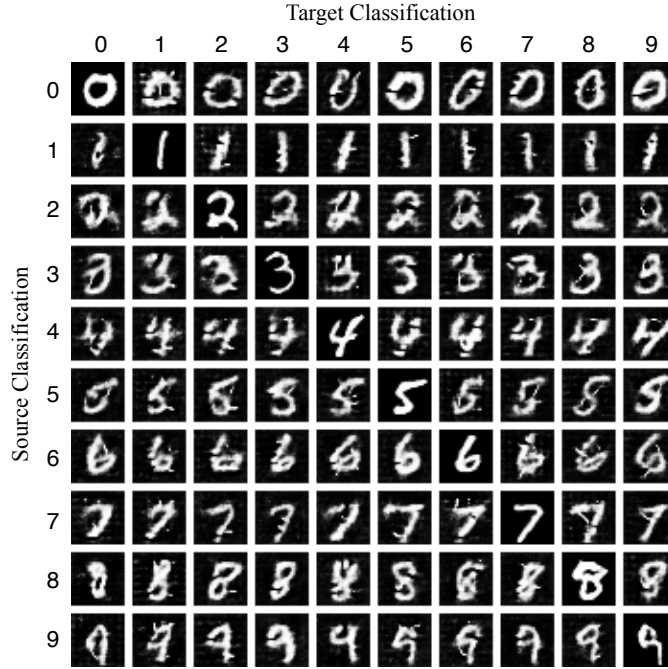


Figure 5: Adversarial examples generated by AT-GAN to different targets on MNIST with the same random noise input for each row. The images on the diagonal are generated by $G_{original}$ which are not adversarial examples and are treated as the initial instances for AT-GAN.

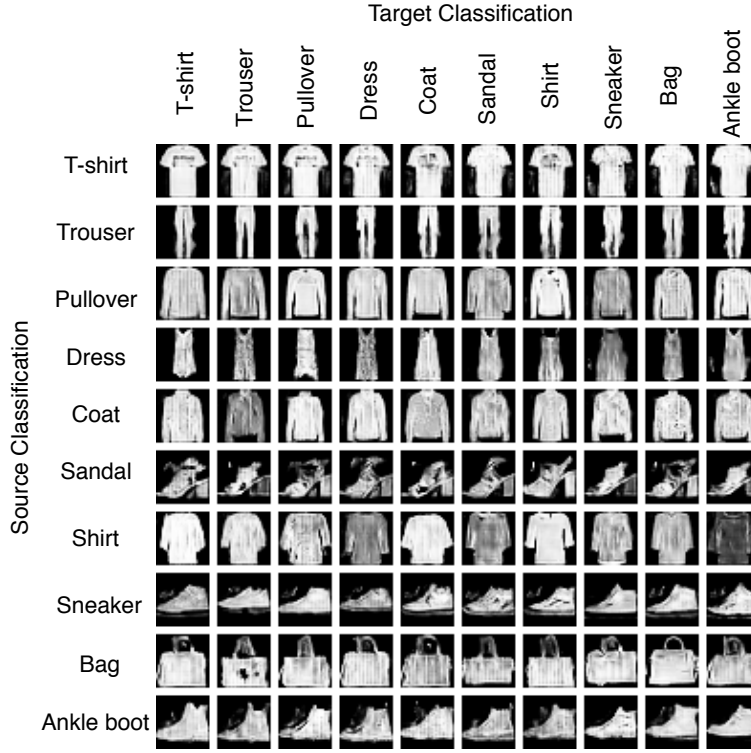


Figure 6: Adversarial examples generated by AT-GAN to different targets on Fashion-MNIST with the same random noise input for each row. The images on the diagonal are generated by $G_{original}$ which are not adversarial examples and are treated as the initial instances for AT-GAN.