# Improving Adversarial Robustness Through the Contrastive-Guided Diffusion Process

Yidong Ouyang <sup>1</sup> Liyan Xie <sup>1</sup> Guang Cheng <sup>2</sup>

#### **Abstract**

Synthetic data generation has become an emerging tool to help improve the adversarial robustness in classification tasks, since robust learning requires a significantly larger amount of training samples compared with standard classification. Among various deep generative models, the diffusion model has been shown to produce high-quality synthetic images and has achieved good performance in improving the adversarial robustness. However, diffusion-type methods are generally slower in data generation as compared with other generative models. Although different acceleration techniques have been proposed recently, it is also of great importance to study how to improve the sample efficiency of synthetic data for the downstream task. In this paper, we first analyze the optimality condition of synthetic distribution for achieving improved robust accuracy. We show that enhancing the distinguishability among the generated data is critical for improving adversarial robustness. Thus, we propose the Contrastive-Guided Diffusion Process (Contrastive-DP), which incorporates the contrastive loss to guide the diffusion model in data generation. We validate our theoretical results using simulations and demonstrate the good performance of Contrastive-DP on image datasets.

#### 1. Introduction

The success of most deep learning methods relies heavily on a massive amount of training data, which can be expensive to acquire in practice. For example, in applications like autonomous driving (O'Kelly et al., 2018) and medical

Proceedings of the 40<sup>th</sup> International Conference on Machine Learning, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

diagnosis (Das et al., 2022), the number of rare scenes is usually very limited in the training dataset. Moreover, the number of labeled data for supervised learning could also be limited in some applications since it may be expensive to label the data. These challenges call for methods that can produce additional data that are easy to generate and can help improve downstream task performance. Synthetic data generation based on deep generative models has shown promising performance recently to tackle these challenges (Sehwag et al., 2022; Gowal et al., 2021; Das et al., 2022).

In synthetic data generation, one aims to learn a synthetic distribution (from which we generate synthetic data) that is close to the true date-generating distribution and, most importantly, can help improve the downstream task performance. Synthetic data generation is highly related to generative models. Among various kinds of generative models, the score-based model and diffusion type models have gained much success in image generation recently (Song & Ermon, 2019; Song et al., 2021b; 2020; Song & Ermon, 2020; Sohl-Dickstein et al., 2015; Nichol & Dhariwal, 2021; Bao et al., 2022; Rombach et al., 2022; Nie et al., 2022; Sun et al., 2022). As validated in image datasets, the prototype of diffusion models, the Denoising Diffusion Probabilistic Model (DDPM) (Ho et al., 2020), and many variants can generate high-quality images as compared with classical generative models such as generative adversarial networks (Dhariwal & Nichol, 2021).

This paper mainly focuses on the adversarial robust classification of image data, which typically requires more training data than standard classification tasks (Carmon et al., 2019). In (Gowal et al., 2021), 100M high-quality synthetic images are generated by DDPM and achieve the state-of-the-art performance on adversarial robustness on the CIFAR-10 dataset, which demonstrates the effectiveness of diffusion models in improving adversarial robustness. However, a major drawback of diffusion-type methods is the slow computational speed. More specifically, DDPM is usually 1000 times slower than GAN (Song et al., 2021a), and this drawback is more serious when generating a large number of samples, e.g., it takes more than 99 GPU days <sup>1</sup> for generating 100M image data according to (Gowal et al., 2021).

<sup>&</sup>lt;sup>1</sup>School of Data Science, The Chinese University of Hong Kong, Shenzhen, China <sup>2</sup>Department of Statistics, University of California, Los Angeles, USA. Correspondence to: Liyan Xie <xieliyan@cuhk.edu.cn>.

<sup>&</sup>lt;sup>1</sup>Running on a RTX 4x2080Ti GPU cluster.

Moreover, the computational cost will increase dramatically when the resolution of images increases, which inspires a plentiful of works studying how to accelerate the diffusion models (Song et al., 2021a; Watson et al., 2022; Ma et al., 2022; Salimans & Ho, 2022; Bao et al., 2022; Cao et al., 2022; Yang et al., 2022). In this paper, we aim to study the aforementioned problem from a different perspective – "how to generate effective synthetic data that are most helpful for the downstream task?" We analyze the *optimal synthetic distribution* for the downstream tasks to improve the sample efficiency of the generative model.

We first study the theoretical insights for finding the optimal synthetic distributions for achieving adversarial robustness. Following the setting considered in (Carmon et al., 2019), we introduce a family of synthetic distributions controlled by the distinguishability of the representation from different classes. Our theoretical results show that the more distinguishable the representation is for the synthetic data, the higher the classification accuracy we will get. Motivated by the theoretical insights, we propose the Contrastive-Guided Diffusion Process (Contrastive-DP) for efficient synthetic data generation, incorporating the gradient of the contrastive learning loss (van den Oord et al., 2018; Chuang et al., 2020; Robinson et al., 2021) into the diffusion process. We conduct comprehensive simulations and experiments on real image datasets to demonstrate the effectiveness of the proposed Contrastive-DP method.

The remainder of the paper is organized as follows. Section 2 presents the problem formulation and preliminaries on diffusion models. Section 3 contains the theoretical insights of optimal synthetic distribution under the Gaussian setting. Motivated by the theoretical insights, Section 4 proposes a new type of synthetic data generation procedure that combines contrastive learning with diffusion models. Finally, Section 5 conducts extensive numerical experiments to validate the good performance of the proposed generation method on simulation and image datasets.

# 2. Problem Formulation and Preliminaries

We first give a brief overview of adversarial robust classification, which is our main focus in this work. It is worth mentioning that the whole framework can be applied to other downstream tasks in general. Denote the feature space as  $\mathcal{X}$ , the corresponding label space as  $\mathcal{Y}$ , and the true (joint) data distribution as  $\mathcal{D} = \mathcal{D}_{\mathcal{X} \times \mathcal{Y}}$ .

Assume we have labeled training data  $\mathcal{D}_{\text{train}} := \{(x_i, y_i)\}_{i=1}^n \text{ sampled from } \mathcal{D}.$  We aim to learn a robust classifier  $f_{\theta} : \mathcal{X} \mapsto \mathcal{Y}$ , parameterized by a learnable  $\theta$ , that can achieve the minimum adversarial loss:

$$\min_{\boldsymbol{\theta}} \mathcal{L}_{adv}(\boldsymbol{\theta}) := \mathbb{E}_{(\boldsymbol{x},y) \sim \mathcal{D}} \Big( \max_{\boldsymbol{\delta} \in \Lambda} \ell(\boldsymbol{x} + \boldsymbol{\delta}, y, \boldsymbol{\theta}) \Big), \quad (1)$$

where  $\ell(x,y,\theta)=\mathbf{1}\{y\neq f_{\theta}(x)\}$  is the 0-1 loss function,  $\mathbf{1}\{\cdot\}$  is the indicator function, and  $\Delta=\{\delta:\|\delta\|_{\infty}\leq\epsilon\}$  is the adversarial set defined using  $\ell_{\infty}$ -norm. Intuitively, the solution to problem (1) is a robust classifier that minimizes the worst-case loss within an  $\epsilon$ -neighborhood of the input features.

In the canonical form of adversarial training, we train the robust classifier  $f_{\theta}$  on the training set  $\mathcal{D}_{\text{train}} := \{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$  by solving the following sample average approximation of the population loss in (1):

$$\min_{\boldsymbol{\theta}} \widehat{\mathcal{L}}_{adv}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^{n} \max_{\boldsymbol{\delta}_i \in \Delta} \ell(\boldsymbol{x}_i + \boldsymbol{\delta}_i, y_i, \boldsymbol{\theta}). \quad (2)$$

### 2.1. Adversarial Training Using Synthetic Data

As shown in (Carmon et al., 2019), adversarial training requires more training data in order to achieve the desired accuracy. Synthetic data generation has been used as a method to artificially increase the size of the training set by generating a sufficient amount of additional data, thus helping improve the learning algorithm's performance (Gowal et al., 2021).

The mainstream generation procedures can be categorized into two types: unconditional and conditional generation. In the unconditional generation, we first generate the features (x) and then assign pseudo labels to them. In the conditional generation, we generate the features conditioned on the desired label. Our analysis is mainly based on the former paradigm, which can be easily generalized to the conditional generation procedure, and our proposed algorithm is also flexible enough for both pipelines.

Denote the distribution of the generated features as  $\widetilde{\mathcal{D}}_{\mathcal{X}}$  and the generated synthetic data as  $\mathcal{D}_{\text{syn}} := \{(\tilde{x}_i, \tilde{y}_i)\}_{i=1}^{\tilde{n}}$ . Here the feature values  $\tilde{x}_i$  are generated from the synthetic distribution  $\widetilde{\mathcal{D}}_{\mathcal{X}}$ , and  $\tilde{y}_i$  are pseudo labels assigned by a classifier learned on the training data  $\mathcal{D}_{\text{train}}$ . Combining the synthetic and real data, we will learn the robust classifier using a larger training set  $\mathcal{D}_{\text{all}} := \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{syn}}$  which now contains  $n + \tilde{n}$  samples:

$$\min_{\boldsymbol{\theta}} \left\{ \eta \left( \frac{1}{n} \sum_{i=1}^{n} \max_{\boldsymbol{\delta}_{i} \in \Delta} \ell(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}, y_{i}, \boldsymbol{\theta}) \right) + (1 - \eta) \left( \frac{1}{\tilde{n}} \sum_{i=1}^{n} \max_{\boldsymbol{\delta}_{i} \in \Delta} \ell(\tilde{\boldsymbol{x}}_{i} + \boldsymbol{\delta}_{i}, \tilde{y}_{i}, \boldsymbol{\theta}) \right) \right\},$$
(3)

where  $\eta \in (0,1)$  is a parameter controlling the weights of synthetic data.

#### 2.2. Diffusion Model for Data Generation

We build our proposed generation procedure based on the Denoising Diffusion Probabilistic Model (DDPM) (Ho et al., 2020) and its accelerated variant Denoising Diffusion Implicit Model (DDIM) (Song et al., 2021a). In the following, we briefly review the key components of DDPM.

The core of DDPM is a forward Markov chain with Gaussian transitions distributions  $q(\boldsymbol{x}_t|\boldsymbol{x}_{t-1})$  to inject Gaussian noise to the original data distribution  $q(\boldsymbol{x}_0)$ . More specifically, (Ho et al., 2020) model the forward Gaussian transition as:

$$q\left(\boldsymbol{x}_{t}|\boldsymbol{x}_{t-1}\right) := \mathcal{N}\left(\sqrt{\alpha_{t}}\boldsymbol{x}_{t-1},\left(1-\alpha_{t}\right)\mathbb{I}\right),$$

where  $\alpha_t, t = 1, 2, \ldots, T$  is a decreasing sequence to control the variance of injected noise, and  $\mathbb{I}$  is the identity covariance matrix. The joint likelihood for the above Markov chain can be written as  $q\left(\boldsymbol{x}_{0:T}\right) = q\left(\boldsymbol{x}_{0}\right)\prod_{t=1}^{T}q\left(\boldsymbol{x}_{t}|\boldsymbol{x}_{t-1}\right)$ . DDPM then propose to use  $p_{\theta}\left(\boldsymbol{x}_{0:T}\right) = p_{\theta}\left(\boldsymbol{x}_{T}\right)\prod_{t=1}^{T}p_{\theta}\left(\boldsymbol{x}_{t-1}|\boldsymbol{x}_{t}\right)$  to model the reverse process, where  $p_{\theta}(\boldsymbol{x}_{t-1}|\boldsymbol{x}_{t})$  is parameterized using a neural network. The training objective is to minimize the Kullback–Leibler (KL) divergence between the forward and reverse process:  $D_{\text{KL}}(q\left(\boldsymbol{x}_{0:T}\right),p_{\theta}\left(\boldsymbol{x}_{0:T}\right))$ , which can be simplified as:

$$\min_{\theta} \mathbb{E}_{t,\boldsymbol{x}_{0},\epsilon} \left[ \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} \left( \sqrt{\bar{\alpha}_{t}} \boldsymbol{x}_{0} + \sqrt{1 - \bar{\alpha}_{t}} \boldsymbol{\epsilon}, t \right) \right\|^{2} \right],$$

where the expectation is taken with respect to  $x_0 \sim q(x_0)$ ,  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$ , and t uniformly distributed in  $\{1, \ldots, T\}$ . Here  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$  and  $\epsilon_{\theta}(x, t)$  denotes the neural network parameterized by  $\theta$ . We refer to (Ho et al., 2020) for the detailed derivation and learning algorithms.

After learning the time-reversed process parameterized by  $\theta$ , the original generation process in (Ho et al., 2020) is a time-reversed Markov chain as follows:

$$\boldsymbol{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \boldsymbol{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta} \left( \boldsymbol{x}_t, t \right) \right) + \sigma_t \boldsymbol{z}_t,$$

starting from  $x_T \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$  and calculating for  $t = T, T - 1, \ldots, 1$ . The output value  $x_0$  is the generated synthetic data. Here  $z_t \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$  if t > 1 and  $z_t = \mathbf{0}$  if t = 1. DDIM (Song et al., 2021a) speeds up the above procedure by generalizing the diffusion process to a non-Markovian process, leading to a sampling trajectory much shorter than T. DDIM carefully designs the forward transition such that  $q(x_t|x_0) = \mathcal{N}\left(\sqrt{\alpha_t}x_0, (1-\alpha_t)\,\mathbb{I}\right)$  for all  $t = 1, \ldots, T$ . The great advantage of DDIM is that it admits the same training objective as DDPM, which means we can adapt the pre-trained model of DDPM and accelerate the sampling process without additional cost. The key sample-generating step in DDIM is as follows:

$$x_{t-1} = \sqrt{\alpha_{t-1}} \underbrace{\left(\frac{x_t - \sqrt{1 - \alpha_t} \epsilon_{\theta} (x_t, t)}{\sqrt{\alpha_t}}\right)}_{\text{predicted } x_0} + \underbrace{\sqrt{1 - \alpha_{t-1}} \cdot \epsilon_{\theta} (x_t, t)}_{\text{pointing to } x_t},$$

$$(4)$$

in which we can generate  $x_{t-1}$  using  $x_t$  and  $x_0$ . Also, the generating process becomes deterministic.

# 3. Theoretical Insights: Optimal Synthetic Distribution

In this section, we consider a concrete distributional model as used in (Carmon et al., 2019; Schmidt et al., 2018), and demonstrate the advantage of refining the synthetic data generation process – using the optimal distribution for synthetic data generation can help reduce the sample complexity needed for robust classification. This provides theoretical insights and motivates the proposed Contrastive-DP method to be introduced in Section 4.

#### 3.1. Theoretical Setup

Consider a binary classification task where  $\mathcal{X} = \mathbb{R}^d$ ,  $\mathcal{Y} = \{-1,1\}$ . The true data distribution  $\mathcal{D}$  is specified as follows. The marginal distribution for label y is uniform in  $\mathcal{Y}$ , and the conditional distribution of features is  $\mathbf{x}|y \sim \mathcal{N}(y\boldsymbol{\mu},\sigma^2\mathbb{I}_d)$ , where  $\boldsymbol{\mu} \in \mathbb{R}^d$  is non-zero, and  $\mathbb{I}_d$  is the d dimensional identity covariance matrix. Thus the marginal feature distribution  $\mathcal{D}_{\mathcal{X}}$  is a Gaussian mixture, for convenience we denote as  $0.5\mathcal{N}(\boldsymbol{\mu},\sigma^2\mathbb{I}) + 0.5\mathcal{N}(-\boldsymbol{\mu},\sigma^2\mathbb{I})$ . Suppose we also generate a set of synthetic data from another synthetic distribution  $\widetilde{\mathcal{D}}$  which could be different from  $\mathcal{D}$ .

We focus on learning a robust linear classifier under the above setting. The family of linear classifiers is represented as  $f_{\theta}(\boldsymbol{x}) = \operatorname{sign}(\boldsymbol{\theta}^{\top}\boldsymbol{x})$ . Recall that we first generate features and then assign pseudo labels to the features. Therefore, a self-learning paradigm is adopted here (Wei et al., 2020). Given a set of unlabeled synthetic features  $\{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_{\tilde{n}}\}$ , we apply an intermediate linear classifier parameterized by  $\hat{\theta}_{\text{inter}} = \frac{1}{n} \sum_{i=1}^{n} y_i x_i$ , learned from real data  $\mathcal{D}_{\text{train}}$ , to assign the pseudo-label. Then, the synthetic data  $\mathcal{D}_{\text{syn}} = \{(\tilde{x}_1, \tilde{y}_1), \dots, (\tilde{x}_{\tilde{n}}, \tilde{y}_{\tilde{n}})\}$ , where  $\tilde{y}_i = \operatorname{sign}(\hat{\theta}_{\text{inter}}^{\top}x_i)$ ,  $i = 1, \dots, \tilde{n}$ . We combine the real data and synthetic data  $\mathcal{D}_{\text{all}} := \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{syn}} = \{(x_i, y_i)\}_{i=1}^n, \{(\tilde{x}_i, \tilde{y}_i)\}_{i=1}^{\tilde{n}}\}$  to obtain an approximate optimal solution  $\hat{\theta}_{\text{final}}$  as (Carmon et al., 2019):

$$\hat{\boldsymbol{\theta}}_{\text{final}} = \frac{1}{n+\tilde{n}} \left( \sum_{i=1}^{n} y_i \boldsymbol{x}_i + \sum_{j=1}^{\tilde{n}} \tilde{y}_j \tilde{\boldsymbol{x}}_j \right). \tag{5}$$

Note that the final linear classifier  $\theta_{\text{final}}$  depends on the synthetic data generated from  $\widetilde{\mathcal{D}}$ . We aim to study which synthetic distribution  $\widetilde{\mathcal{D}}$  can help reduce the adversarial classification error (also called robust error)

$$\begin{split} & \text{err}_{\text{robust}}\left(f_{\hat{\boldsymbol{\theta}}_{\text{final}}}\right) := \mathbb{P}_{(\boldsymbol{x},y) \sim \mathcal{D}}(\exists \boldsymbol{\delta} \in \Delta, f_{\hat{\boldsymbol{\theta}}_{\text{final}}}\left(\boldsymbol{x} + \boldsymbol{\delta}\right) \neq y), \\ & \text{where } \Delta = \{\boldsymbol{\delta} : \|\boldsymbol{\delta}\|_{\infty} \leq \epsilon\}. \text{ And we similarly define the standard error as } & \text{err}_{\text{standard}}\left(f_{\hat{\boldsymbol{\theta}}_{\text{final}}}\right) := \mathbb{P}_{(\boldsymbol{x},y) \sim \mathcal{D}}(f_{\hat{\boldsymbol{\theta}}_{\text{final}}}(\boldsymbol{x}) \neq y) \\ & \text{which will be used later.} \end{split}$$

Remark 3.1 (Comparison with existing literature). In (Carmon et al., 2019; Deng et al., 2021), sample complexity results are analyzed based on the same Gaussian setting. The major difference is that they all assume the learned linear classifier  $\hat{\theta}_{\text{final}}$  is only learned from synthetic data  $\mathcal{D}_{\text{syn}}$  rather than the combination of the real and synthetic data  $\mathcal{D}_{\text{all}}$ . In general, our theoretical setup matches well with the practical algorithms.

# 3.2. Theoretical Insights for Optimal Synthetic Distribution

We first study the desired properties of the synthetic distribution  $\widetilde{\mathcal{D}}$  that can lead to a better adversarial classification accuracy when the additional synthetic sample  $\mathcal{D}_{syn}$  is used in the training stage. In (Carmon et al., 2019), the standard case  $\widetilde{\mathcal{D}} = \mathcal{D}$  is studied, i.e., they consider the case that additional unlabeled data from the true distribution  $\mathcal{D}$  is available, and they characterize the usefulness of those additional training data. Compared with (Carmon et al., 2019), we consider general distributions  $\widetilde{\mathcal{D}}$  which does not necessarily equal to  $\mathcal{D}$ .

First note that by the Bayes rule, the optimal decision boundary for the true data distribution is given by  $\boldsymbol{\mu}^{\top}\boldsymbol{x}=0$ . Therefore, we restrict our attention to synthetic data distributions that satisfy: (i) the marginal distribution of the label  $\tilde{y}$  is also uniform in  $\mathcal{Y}$ , same as  $\mathcal{D}$ ; (ii) the conditional probability densities  $p(\tilde{\boldsymbol{x}}|\tilde{y}=1)$  and  $p(\tilde{\boldsymbol{x}}|\tilde{y}=-1)$  of the synthetic data distribution are symmetric around the true optimal decision boundary  $\boldsymbol{\mu}^{\top}\boldsymbol{x}=0$ . More specifically, we start with a special case of the synthetic data distribution  $\tilde{\mathcal{D}}_{\mathcal{X}}=0.5\mathcal{N}(\tilde{\boldsymbol{\mu}},\sigma^2\mathbb{I})+0.5\mathcal{N}(-\tilde{\boldsymbol{\mu}},\sigma^2\mathbb{I})$  (note that when  $\tilde{\boldsymbol{\mu}}=c\boldsymbol{\mu}$  for some constant c, the above two conditions are all satisfied).

In the following proposition, we present several representative scenarios of synthetic distributions in terms of how they may contribute to the downstream classification task. Figure 1 gives a pictorial demonstration for different cases.

**Proposition 3.2.** Consider a special form of synthetic distributions  $\widetilde{\mathcal{D}}_{\mathcal{X}} = 0.5 \mathcal{N}(\widetilde{\boldsymbol{\mu}}, \sigma^2 \mathbb{I}) + 0.5 \mathcal{N}(-\widetilde{\boldsymbol{\mu}}, \sigma^2 \mathbb{I})$  and assume  $\{\widetilde{\boldsymbol{x}}_1, \dots, \widetilde{\boldsymbol{x}}_{\widetilde{n}}\}$  are samples from  $\widetilde{\mathcal{D}}_{\mathcal{X}}$ . We follow the self-learning paradigm described in Section 3.1 to learn the classifier  $f_{\widehat{\boldsymbol{\theta}}_{nool}}$ , when  $\widetilde{n}$  is sufficiently large we have:

Case 1: Inefficient  $\widetilde{\mathcal{D}}_{\mathcal{X}}$ . When  $\langle \widetilde{\boldsymbol{\mu}}, \boldsymbol{\mu} \rangle = 0$ , the standard error  $\operatorname{err}_{standard}(f_{\widehat{\boldsymbol{\theta}}_{final}})$  achieves the maximum and when  $\langle \widetilde{\boldsymbol{\mu}}, \boldsymbol{\mu} - \varepsilon \mathbf{1}_d \rangle = 0$ , the robust error  $\operatorname{err}_{robust}(f_{\widehat{\boldsymbol{\theta}}_{final}})$  achieves the maximum.

Case 2: Optimal  $\widetilde{\mathcal{D}}_{\mathcal{X}}$  for clean accuracy. When  $\widetilde{\boldsymbol{\mu}} = c\boldsymbol{\mu}$  for c > 0,  $\operatorname{err}_{standard}(f_{\widehat{\boldsymbol{\theta}}_{final}})$  achieves the minimum, and the larger the c is, the smaller the  $\operatorname{err}_{standard}(f_{\widehat{\boldsymbol{\theta}}_{final}})$ .

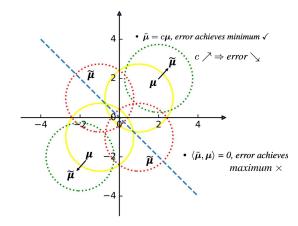


Figure 1. Demonstration of Proposition 3.2. We refer to the main text for a detailed explanation.

Case 3: Optimal  $\widetilde{\mathcal{D}}_{\mathcal{X}}$  for robust accuracy. When  $\widetilde{\boldsymbol{\mu}} = c(\boldsymbol{\mu} - \varepsilon \mathbf{1}_d)$  for c > 0, the robust error  $\operatorname{err}_{robust}(f_{\widehat{\boldsymbol{\theta}}_{final}})$  achieves the minimum, and the larger the c is, the smaller the  $\operatorname{err}_{robust}(f_{\widehat{\boldsymbol{\theta}}_{final}})$ .

Remark 3.3 (Comparison with the existing characterization of the synthetic distribution). We briefly comment on the main differences and similarities with (Deng et al., 2021), in which a similar result was presented in Theorem 4 therein. In (Deng et al., 2021), the final solution of  $\theta^*$  was given for minimizing robust error  $\text{err}_{\text{robust}}\left(\hat{f}_{\hat{\theta}_{\text{final}}}\right)$  and they provides a specific unlabeled distribution  $\tilde{\mu} = \mu - \varepsilon \mathbf{1}_d$  that achieves better performance under certain condition. In this paper, we propose a general family of optimal distribution controlled by a scalar c, which represents the distinguishability of the feature. The final  $\theta^*$  used in (Deng et al., 2021) can be viewed as a special case of c=1. Therefore, our conclusion points out the optimality condition of unlabeled distribution and inspires a line of work to improve the performance of  $\hat{\theta}_{\text{final}}$  by making the feature of unlabeled distribution distinguishable.

We also study the sample complexity for the synthetic distributions satisfying the condition in Proposition 3.2. The results below show that for larger c, we typically need fewer synthetic samples to achieve the desired robust accuracy.

**Theorem 3.4.** Under the parameter setting  $\epsilon < \frac{1}{2}$ ,  $\sigma = (nd)^{1/4}$ ,  $\|\mu\|^2 = d$ , there exists a universal constant  $\tilde{C}$  such that for  $\epsilon^2 \sqrt{d/n} \geq \tilde{C}$ , where n is the number of labeled real data used to construct the intermediate classifier, and additional  $\tilde{n}$  synthetic feature generated with mean vector  $\pm c\mu$  and pseudo labels, if

$$\tilde{n} \ge \frac{288n}{c} \epsilon^2 \sqrt{\frac{d}{n}},$$

then

$$\mathbb{E}_{\hat{\boldsymbol{\theta}}_{\mathit{final}}}\left[\;\mathit{err}_{\mathit{robust}}\left(f_{\hat{\boldsymbol{\theta}}_{\mathit{final}}}\right)\right] \leq 10^{-3}.$$

Table 1. Simulation results validating findings in Proposition 3.2 and Theorem 3.4. We use "Real" to denote the real data distribution and n to denote the number of data from the real distribution, and we use "c" to denote different synthetic distributions with mean  $\tilde{\mu}=c\mu$  and use  $\tilde{n}$  to denote the number of synthetic data. The average accuracy and the standard deviation in the bracket are obtained from 50 repetitions.

		clean acc	rob acc
Real	$\begin{array}{ c c } n = 10 \\ n = 100 \end{array}$	0.9023 (0.0192) 0.9682 (0.0014)	0.6843 (0.0359) 0.8239 (0.0028)
c = 0.5		0.7562 (0.0564) 0.9505 (0.0047)	0.4611 (0.0694) 0.7848 (0.0111)
c = 1		0.8866 (0.0273) 0.9695 (0.0012)	0.6557 (0.0487) 0.8239 (0.0031)
c = 1.5		<b>0.9400</b> (0.0100) <b>0.9743</b> (0.0008)	<b>0.7603</b> (0.0233) <b>0.8343</b> (0.0011)

**Simulation Results.** To verify the findings in Proposition 3.2 and Theorem 3.4, we conduct extensive simulation experiments under Gaussian distributions with varying data dimensions, sample sizes, and the mean vector  $\tilde{\mu}$ . In most cases, we find increasing c for synthetic distribution can lead to better clean and robust accuracy. A detailed description of the experimental setting can be found in Appendix B. In Table 1, we demonstrate the clean and robust accuracy learned on synthetic distribution with the angle between  $\mu$  and  $\epsilon \mathbf{1}_d$  equals  $0^{\circ}$ . Remarkably, the classifier learned only from the synthetic distribution with  $\tilde{\mu} = c\mu$  with c > 1 achieves better performance even than the iid samples (denoted as "Real" in Table).

The closed form of optimal synthetic distribution is  $\tilde{\mu} = \mu - \varepsilon \mathbf{1}_d$  as stated in Proposition 3.2. It is interesting to see when  $\tilde{\mu} = c\mu$  and  $\mu$  is not aligned on  $\epsilon \mathbf{1}_d$ , whether increasing c provides better data distribution for learning a robust classifier. The answer is still true, demonstrated by the experiment results (with  $30^\circ$ ,  $60^\circ$ , and  $90^\circ$ ) in Table 4, 5, 6 and 7 in Appendix B.

#### 4. Contrastive-Guided Diffusion Process

It has been shown in Proposition 3.2 and Theorem 3.4 that the synthetic data can help improve the classification task, especially when the representation of different classes is more distinguishable in the synthetic distribution. Meanwhile, the contrastive loss (van den Oord et al., 2018) can be adopted to explicitly control the distances of the representation of different classes. Therefore, we propose a variant of the classical diffusion model, named *Contrastive-Guided Diffusion Process (Contrastive-DP)*, to enhance the sample efficiency of the generative model. In this section, we first present the overall algorithm of the proposed Contrastive-

**Algorithm 1** Generation in Contrastive-guided Diffusion Process (Contrastive-DP)

**Require:** Contrastive loss temperature  $\tau$ , diffusion process

hyperparameter  $\sigma_t$ 1:  $\mathbb{X}_T = \{\boldsymbol{x}_T^{(i)}\}_{i=1}^m \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$ 2: t = T3: while  $t \neq 1$  do

4: for i = 1 to m do

5: Sampling  $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$ 6: Choosing  $\boldsymbol{x}_p^{(i)}$  as the positive pair of  $\boldsymbol{x}_t^{(i)}$ 7:  $\Delta \boldsymbol{x}_t^{(i)} = \lambda \cdot \nabla_{\boldsymbol{x}_t^{(i)}} \ell_{\text{contra}}(\boldsymbol{x}_t^{(i)}, \boldsymbol{x}_p^{(i)}; \tau) + \boldsymbol{\epsilon}_{\theta}(\boldsymbol{x}_t^{(i)}, t)$ 8:  $\boldsymbol{x}_{t-1}^{(i)} = \sqrt{\alpha_{t-1}} (\frac{\boldsymbol{x}_t^{(i)} - \sqrt{1 - \alpha_t \Delta \boldsymbol{x}_t^{(i)}}}{\sqrt{\alpha_t}}) + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \Delta \boldsymbol{x}_t^{(i)} + \sigma_t \boldsymbol{\epsilon}_t$ 9: t = t - 1

DP procedure in Section 4.1, then we describe the detailed design of the contrastive loss in Section 4.2.

#### 4.1. Algorithm for Contrastive-Guided DP

end for

12: Return  $X_0 = \{x_0^{(i)}\}_{i=1}^m$ 

11: end while

10:

The detailed generation procedure of Contrastive-DP is given in Algorithm 1. We highlight below some major differences between the proposed Contrastive-DP and the vanilla DDIM algorithm. In each time step t of the generation procedure, given the current value  $\boldsymbol{x}_t^{(i)}$ , we add the gradient of the contrastive loss  $\ell_{\text{contra}}(\boldsymbol{x}_t^{(i)}, \boldsymbol{x}_p^{(i)}; \tau)$  with respect to  $\boldsymbol{x}_t^{(i)}$  to the original diffusion generative process, here  $\boldsymbol{x}_p^{(i)}$  is the positive pair of  $\boldsymbol{x}_t^{(i)}$  (will be explained in detail later),  $\tau$  is the temperature for softmax, and  $\lambda$  is the hyperparameter balancing the contrastive loss within the diffusion process.

This modification ensures that the generated data will be distinguishable among data in the same batch. The construction of the contrastive loss  $\ell_{contra}(\cdot)$  is very flexible – we can adopt multiple forms of contrastive loss together with different selection strategies of positive and negative pairs, which will be discussed in detail in the following.

#### 4.2. Contrastive Loss for Diffusion Process

Let  $\mathbb{X} = \{ \boldsymbol{x}_1, ..., \boldsymbol{x}_m \}$  be a minibatch of training data. We apply the contrastive loss to the embedding space. Assume  $f(\cdot)$  is the feature extractor that maps the input data in  $\mathbb{X}$  onto the embedding space. In general, we adopt two forms of the contrastive loss  $\ell_{\text{contra}}(\boldsymbol{x}_t^{(i)}, \boldsymbol{x}_p^{(i)}; \tau)$  which will be used in Algorithm 1.

First is the InfoNCE loss:

$$\ell_{\text{InfoNCE}}\left(\boldsymbol{x}_{a}, \boldsymbol{x}_{p}; \tau\right) = -\log(\frac{g_{\tau}(\boldsymbol{x}_{a}, \boldsymbol{x}_{p})}{\sum_{k=1}^{m} \mathbf{1}_{k \neq a} g_{\tau}(\boldsymbol{x}_{a}, \boldsymbol{x}_{k})}),$$

where m is the batch size,  $\tau$  is the temperature for softmax,  $\boldsymbol{x}_a, \boldsymbol{x}_p$  denote the anchor and the positive pair, respectively,  $g_{\tau}(\boldsymbol{x}, \boldsymbol{x}') = \exp(f(\boldsymbol{x})^{\top} f(\boldsymbol{x}')/\tau)$ , and all images except the anchor  $\boldsymbol{x}_a$  in the minibatch  $\mathbb{X}$  is negative pairs. InfoNCE loss is an unsupervised learning metric and does not explicitly distinguish the representation from different classes, which implicitly regards the representation from the same class as negative pair.

Second is the hard negative mining loss:

$$\ell_{\text{HNM}}\left(\boldsymbol{x}_{a}, \boldsymbol{x}_{p}; \tau\right) = -\log \frac{g_{\tau}(\boldsymbol{x}_{a}, \boldsymbol{x}_{p})}{g_{\tau}(\boldsymbol{x}_{a}, \boldsymbol{x}_{p}) + \frac{m}{\tau^{-}} h_{\tau}(\boldsymbol{x}_{a})},$$

where

$$h_{\tau}(\boldsymbol{x}_a) = \mathbb{E}_{\boldsymbol{x}_n \sim q_{\beta}}[g_{\tau}(\boldsymbol{x}_a, \boldsymbol{x}_n)] - \tau^{+} \mathbb{E}_{\boldsymbol{v} \sim q_{\beta}^{+}}[g_{\tau}(\boldsymbol{x}_a, \boldsymbol{v})],$$

and m denotes the batch size,  $\tau^-=1-\tau^+$  denotes the probability of observing any different class with  $x_a$  and  $q_\beta$  is an unnormalized von Mises–Fisher distribution (Jammalamadaka, 2011), with mean direction f(x) and "concentration parameter"  $\beta$  to control the hardness of negative mining;  $q_\beta$  and  $q_\beta^+$  can be easily approximated by Monte-Carlo importance sampling techniques. We refer to (Chuang et al., 2020; Robinson et al., 2021) for detailed descriptions of hard negative mining contrastive loss. Compared with the InfoNCE loss that does not consider class/label information, the hard negative mining (HNM) loss enhances the discriminative ability of different classes in the feature space.

It is worth mentioning that the Contrastive-DP enjoys the plugin-type property – it does not modify the original training procedure of diffusion processes and can be easily adopted to various kinds of diffusion models.

Moreover, it has a close relationship with Stein Variational Gradient Descent (Liu & Wang, 2016; Liu, 2017). By utilizing the existing theoretical tools (Shi & Mackey, 2022) for analyzing the convergence rate for SVGD, we may also get the convergence guarantee for Contrastive-DP, which is left to future work.

*Remark* 4.1 (Connection with Stein Variational Gradient Descent). Recall the updating step in Contrastive-DP is

$$\Delta \boldsymbol{x}_t^{(i)} = \underbrace{\boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\boldsymbol{x}_t^{(i)}, t)}_{\text{score}} + \lambda \underbrace{\nabla_{\boldsymbol{x}_t^{(i)}} \ell_{\text{contra}}(\boldsymbol{x}_t^{(i)}, \boldsymbol{x}_p^{(i)}; \tau)}_{\text{repulsive force}}.$$

And the updating step  $\Delta x_t^{(i)}$  in SVGD equals (Liu, 2017)

$$\frac{1}{n} \sum_{j=1}^{n} \bigg[ \underbrace{s_p \big( \boldsymbol{x}_t^{(j)} \big) k \big( \boldsymbol{x}_t^{(j)}, \boldsymbol{x}_t^{(i)} \big)}_{\text{weighted sum of score}} + \underbrace{\nabla_{\boldsymbol{x}_t^{(j)}} k \big( \boldsymbol{x}_t^{(j)}, \boldsymbol{x}_t^{(i)} \big)}_{\text{repulsive force}} \bigg],$$

where  $\boldsymbol{x}_t^{(j)}$  is other samples in the batch,  $s_p(\boldsymbol{x}) := \nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})$  is the score function, and  $k(\boldsymbol{x}, \boldsymbol{x}')$  is a positive definite kernel. Contrastive-DP is similar to the SVGD with the score that pulls the particles to high-density region and the repulsive force that pulls the particles away from each other, but the score in Contrastive-DP is easier to calculate as it does not require a weighted sum over the current batch.

Numerical Validations. We first demonstrate the effectiveness of Contrastive-DP in Figure 2 using a simulation example. Consider the binary classification problem as in Section 3.1, and the real data for each class are generated from a Gaussian distribution. Figure 2(a) demonstrates the synthetic data generated by the vanilla diffusion model, which recovers the ground-truth Gaussian distribution well. When using the contrastive-DP procedure with HNM loss, we obtain the generated synthetic data as shown in Figure 2(b), which is more distinguishable with a much smaller variance.

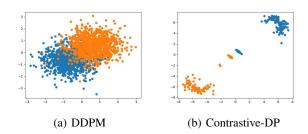


Figure 2. An illustration of the effectiveness of synthetic distribution guided by contrastive loss. More details can be found in Appendix C.2.

In addition, Figure 6 and Figure 7 in Appendix C.2 demonstrate the synthetic data distribution guided by different kinds of contrastive loss mentioned above. It can be shown that InfoNCE loss and hard negative mining method cannot explicitly distinguish the data within the same class and thus form a circle within each class to maximize the distance between samples, while the conditional version of contrastive loss (given the oracle class information) can make two classes more separable.

# 5. Real Data Examples

In this section, we demonstrate the effectiveness of the proposed contrastive guided diffusion process for synthetic data generation in adversarial classification tasks. We first compare the performance of Contrastive-DP with the vanilla diffusion models in Section 5.1. Then, we present a comprehensive ablation study on the performance of Contrastive-DP to shed insights on how to adopt the contrastive loss functions and further data selection methods on the diffu-

sion model in Section 5.2.

#### 5.1. Experimental Results

We test the contrastive-DP algorithm on three image datasets, the MNIST dataset (LeCun et al., 1998), CIFAR-10 dataset (Krizhevsky, 2009), and the Traffic Signs dataset (Houben et al., 2013). A detailed description of the pipeline for generating data and the corresponding hyperparameter can be found in Appendix D.3.

Figure 3 shows the efficacy of synthetic data in terms of improving adversarial robustness on MNIST data. We fix the total number of images for training the classifier as 5K (e.g., when we use 1K real data and 4K synthetic data, to avoid information leakage, we only use these 1K real data to train the diffusion model. The same case for 2K, 3K, and 4K.) <sup>2</sup> Notably, synthetic data improves the clean and robust accuracy even without using any real data, which is consistent with the results in Proposition 3.2 and Theorem 3.4. Moreover, lots of proposed methods (Madry et al., 2017; Tsipras et al., 2018; Zhang et al., 2019) improves the robust accuracy at the sacrifice of clean accuracy, while adding contrastive guidance increase both the clean and the robust accuracy at the same time under the majority of the settings in Figure 3, which shows the potential of Contrastive-DP in achieving a better trade-off between clean and robust accuracy.

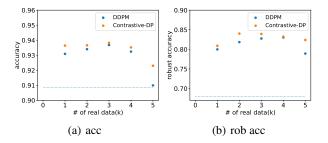


Figure 3. The efficacy of the synthetic data on MNIST dataset. All the results are trained with 5K data (with different proportions of real and synthetic data). Subfigure (a) and (b) show the clean and adversarial accuracy on the MNIST dataset, respectively. The dashed line is the performance training only on 5K real data.

Table 2 demonstrates the effectiveness of our contrastive-DP algorithm on the CIFAR-10 dataset <sup>3</sup>, which achieves better robust accuracy on all data regimes than the vanilla





(a) DDIM

(b) Contrasrive-DP

Figure 4. Comparison of the Contrastive-DP with the vanilla DDIM. The image in the same position on subfigures (a) and (b) has the same initialization.

DDPM and DDIM. All of the results are higher than the baseline result without synthetic data ( $81.98\% \pm 0.58\%$  for clean accuracy and  $50.42\% \pm 0.35\%$  for robust accuracy) by a large margin (+6.18% in 50K setting and +9.57% in 1M setting). Table 3 demonstrates the effectiveness of our contrastive-DP algorithm on the Traffic Signs dataset. Our contrastive-DP achieves better clean and robust accuracy than the vanilla DDPM model and is also higher than the baseline result without synthetic data by a large margin (+10.24%).

To visualize the effectiveness of the guidance, we use the same initialization to generate images by DDIM and Contrastive-DP. We find the guidance of the contrastive loss changes the category of the synthetic images or makes the synthetic images realistic (colorful).

Moreover, we visualize the t-SNE of the finial classifier learned on different synthetic data. We find with the guidance of the contrastive loss, the final classifier learns a better representation that makes the feature of the images from different classes more separable than the final classifier learned on the images generated by the vanilla DDIM.

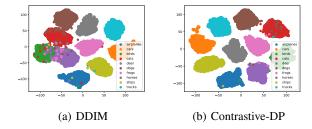


Figure 5. A comparison of the T-SNE of the final classifier learned on different synthetic data on the CIFAR-10 dataset.

#### 5.2. Ablation Studies

In this subsection, we investigate the effectiveness of seven kinds of contrastive loss, the effect of strength of the contrastive loss, and four proposed selection criterion for choosing more informative data from synthetic data. Due to the

<sup>&</sup>lt;sup>2</sup>The case with 5K training images is a special scenario: (i) the accuracy resulted from using all of the 5K real data is the baseline (dashed line shown in Figure 3); (ii) while the accuracy resulted from using 5K synthetic data generated by the diffusion model trained on 5K real data is shown in the scatter plot.

<sup>&</sup>lt;sup>3</sup>Since the Pytorch Implementation of (Gowal et al., 2021) is not open source, we utilize the best unofficial implementation to reconduct all the experiments for a fair comparison.

Table 2. The clean and adversarial accuracy on CIFAR-10 dataset. The robust accuracy is reported by AUTOATTACK (Croce & Hein, 2020) with  $\epsilon_{\infty}=8/255$  and WRN-28-10. 50k, 200k, and 1M denote the number of synthetic used for adversarial training. The results and the standard deviation in the bracket are obtained from 3 repetitions.

	50K		200	0K	1M	
	clean acc	rob acc	clean acc	rob acc	clean acc	rob acc
DDIM	84.05%(0.06%)	56.29%(0.15%)	84.86%(0.43%)	57.83%(0.28%)	85.73%(0.51%)	59.85%(0.26%)
Contrastive-DP	83.66%(0.21%)	<b>56.60%</b> (0.17%)	85.71%(0.18%)	<b>58.24%</b> (0.20%)	<b>86.30</b> %(0.09%)	<b>59.99</b> %(0.23%)
DDPM	84.84%(0.37%)	56.30%(0.06%)	85.23%(0.25%)	58.28%(0.09%)	86.86%(0.04%)	59.03%(0.16%)
Contrastive-DP	84.70%(0.43%)	56.18%(0.10%)	85.61%(0.14%)	<b>58.62%</b> (0.12%)	86.30%(0.10%)	<b>59.74%</b> (0.26%)

*Table 3.* The clean and adversarial accuracy on the Traffic Signs dataset. The results and the standard deviation in the bracket are obtained from 3 repetitions.

	clean acc	rob acc
No additional data	78.52% (0.16%)	46.03% (0.85%)
DDPM	86.79% (0.12%)	56.01% (0.14%)
Contrastive-DP	86.94% (0.32%)	56.27% (0.25%)

space limit, we refer to Appendix E.1, E.2, and E.3 for the detailed numerical results, respectively.

#### 6. Related Work

Using generative models to improve adversarial robustness has attracted increasing attention recently. (Gowal et al., 2021) uses 100M high-quality images generated by DDPM together with the original training set to achieve state-of-the-art performance on the CIFAR-10 dataset. They propose to use Complementary as an important metric for measuring the efficacy of the synthetic data. In (Sehwag et al., 2022), it was claimed that the transferability of adversarial robustness between two data distributions is measured by conditional Wasserstein distance, which inspires us to use it as a criterion for selecting samples. Our work follows the same line, but we investigate how to generate the samples with high information rather than applying the selection to the data generated by the vanilla diffusion model. Below we also summarize some closely related work in different lines.

Sample-Efficient Generation. We can view the sample-efficient generation problem as a Bi-level optimization problem. We can regard how to synthesize data as the meta objective and the performance of the model trained on the synthetic data as the inner objective. For data-augmentation based methods, (Ruiz et al., 2019) adopt a reinforcement learning based method for optimizing the generator in order to maximize the training accuracy. For active learning based methods, (Tran et al., 2019) use an Auto-Encoder to generate new samples based on the informative training data selected by the acquisition function. Besides, (Kim et al., 2020b) combines the active learning criterion with

data augmentation methods. They use the gradient of acquisition function after one-step augmentation as guidance for training the augmentation policy network.

Theoretical Analysis of Adversarial Robustness. In (Schmidt et al., 2018), the sample complexity of adversarial robustness has been shown to be substantially larger than standard classification tasks in the Gaussian setting. (Carmon et al., 2019) bridges this gap by using the self-training paradigm and corresponding unlabeled data. (Deng et al., 2021) further extends the aforementioned conclusion by leveraging out-of-domain unlabeled data. However, there still lacks analysis on the optimal distribution for synthetic data and the corresponding generation algorithm.

Contrastive Learning. Contrastive learning algorithms have been widely used for representation learning (Chen et al., 2020; He et al., 2020; Grill et al., 2020). The vanilla contrastive learning loss, InfoNCE (van den Oord et al., 2018), aims to draw the distance between positive pairs and push the negative pairs away. To mitigate the problem that not all negative pairs may be true negatives, the negative hard mining criterion was proposed in (Chuang et al., 2020; Robinson et al., 2021).

#### 7. Conclusion and Discussion

We delve into which kind of synthetic distribution is optimal for the downstream task, especially for achieving adversarial robustness in image data classification. We derive the optimality condition under the Gaussian setting and propose the Contrastive-guided Diffusion Process (Contrastive-DP), a plug-in algorithm suitable for various types of diffusion models. We verify our theoretical results on the simulated Gaussian example and demonstrate the superiority of the Contrastive-DP algorithm on real image datasets.

It would also be interesting to study the theoretical guarantee of the contrastive-guided diffusion process from the perspective of optimal control. We believe that the proposed plug-in type algorithm can also be generalized to loss functions other than contrastive loss, such as the acquisition function in active learning, for other downstream tasks.

# 8. Acknowledgement

Yidong Ouyang and Liyan Xie are partially supported by UDF01002142 through The Chinese University of Hong Kong, Shenzhen, and Grant J00220220004 through Shenzhen Research Institute of Big Data. Guang Cheng is partially supported by Office of Naval Research, ONR (N00014-22-1-2680), NSF – SCALE MoDL (2134209), and Meta gift fund.

#### References

- Bao, F., Li, C., Zhu, J., and Zhang, B. Analytic-dpm: an analytic estimate of the optimal reverse variance in diffusion probabilistic models. In *ICLR*, 2022.
- Cao, H., Tan, C., Gao, Z., Chen, G., Heng, P.-A., and Li, S. Z. A survey on generative diffusion model. arXiv preprint arXiv:2209.02646, 2022.
- Carmon, Y., Raghunathan, A., Schmidt, L., Liang, P., and Duchi, J. C. Unlabeled data improves adversarial robustness. In *NeurIPS*, 2019.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. E. A simple framework for contrastive learning of visual representations. *ArXiv*, abs/2002.05709, 2020.
- Chuang, C.-Y., Robinson, J., Lin, Y.-C., Torralba, A., and Jegelka, S. Debiased contrastive learning. In *NeurIPS*, 2020.
- Croce, F. and Hein, M. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. *ArXiv*, abs/2003.01690, 2020.
- Das, H. P., Tran, R., Singh, J., Yue, X., Tison, G. H., Sangiovanni-Vincentelli, A. L., and Spanos, C. J. Conditional synthetic data generation for robust machine learning applications with limited pandemic data. In *AAAI*, 2022.
- Deng, Z., Zhang, L., Ghorbani, A., and Zou, J. Y. Improving adversarial robustness via unlabeled out-of-domain data. In *AISTATS*, 2021.
- Dhariwal, P. and Nichol, A. Diffusion models beat gans on image synthesis. *ArXiv*, abs/2105.05233, 2021.
- Fan, L., Liu, S., Chen, P.-Y., Zhang, G., and Gan, C. When does contrastive learning preserve adversarial robustness from pretraining to finetuning? In *Neural Information Processing Systems*, 2021.
- Gowal, S., Qin, C., Uesato, J., Mann, T. A., and Kohli, P. Uncovering the limits of adversarial training against norm-bounded adversarial examples. *ArXiv*, abs/2010.03593, 2020.

- Gowal, S., Rebuffi, S.-A., Wiles, O., Stimberg, F., Calian, D. A., and Mann, T. A. Improving robustness using generated data. In *NeurIPS*, 2021.
- Grill, J.-B., Strub, F., Altch'e, F., Tallec, C., Richemond, P. H., Buchatskaya, E., Doersch, C., Pires, B. Á., Guo, Z. D., Azar, M. G., Piot, B., Kavukcuoglu, K., Munos, R., and Valko, M. Bootstrap your own latent: A new approach to self-supervised learning. *ArXiv*, abs/2006.07733, 2020.
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. B. Momentum contrast for unsupervised visual representation learning. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 9726–9735, 2020.
- Hendrycks, D. and Gimpel, K. Gaussian error linear units (gelus). *arXiv: Learning*, 2016.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. In *NeurIPS*, 2020.
- Houben, S., Stallkamp, J., Salmen, J., Schlipsing, M., and Igel, C. Detection of traffic signs in real-world images:
  The German Traffic Sign Detection Benchmark. In *International Joint Conference on Neural Networks*, number 1288, 2013.
- Izmailov, P., Podoprikhin, D., Garipov, T., Vetrov, D. P., and Wilson, A. G. Averaging weights leads to wider optima and better generalization. *ArXiv*, abs/1803.05407, 2018.
- Jammalamadaka, S. R. Directional statistics, i. 2011.
- Kim, M., Tack, J., and Hwang, S. J. Adversarial self-supervised contrastive learning. *ArXiv*, abs/2006.07589, 2020a.
- Kim, Y.-Y., Song, K., Jang, J., and Moon, I.-C. Lada: Look-ahead data acquisition via augmentation for active learning. *ArXiv*, abs/2011.04194, 2020b.
- Krizhevsky, A. Learning multiple layers of features from tiny images. 2009.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proc. IEEE*, 86:2278–2324, 1998.
- Liu, Q. Stein variational gradient descent as gradient flow. In NIPS, 2017.
- Liu, Q. and Wang, D. Stein variational gradient descent: A general purpose bayesian inference algorithm. In NIPS, 2016.
- Ma, H., Zhang, L., Zhu, X., Zhang, J., and Feng, J. Accelerating score-based generative models for high-resolution image synthesis. *arXiv* preprint arXiv:2206.04029, 2022.

- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. *ArXiv*, abs/1706.06083, 2017.
- Nichol, A. Q. and Dhariwal, P. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pp. 8162–8171. PMLR, 2021.
- Nie, W., Guo, B., Huang, Y., Xiao, C., Vahdat, A., and Anandkumar, A. Diffusion models for adversarial purification. In *International Conference on Machine Learning*, 2022.
- O'Kelly, M., Sinha, A., Namkoong, H., Duchi, J. C., and Tedrake, R. Scalable end-to-end autonomous vehicle testing via rare-event simulation. In *NeurIPS*, 2018.
- Robinson, J., Chuang, C.-Y., Sra, S., and Jegelka, S. Contrastive learning with hard negative samples. In *ICLR*, 2021.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- Ruiz, N., Schulter, S., and Chandraker, M. Learning to simulate. *ArXiv*, abs/1810.02513, 2019.
- Salimans, T. and Ho, J. Progressive distillation for fast sampling of diffusion models. In *ICLR*, 2022.
- Schmidt, L., Santurkar, S., Tsipras, D., Talwar, K., and Madry, A. Adversarially robust generalization requires more data. In *NeurIPS*, 2018.
- Sehwag, V., Mahloujifar, S., Handina, T., Dai, S., Xiang, C., Chiang, M., and Mittal, P. Robust learning meets generative models: Can proxy distributions improve adversarial robustness? In *ICLR*, 2022.
- Shi, J. H. and Mackey, L. W. A finite-particle convergence rate for stein variational gradient descent. *ArXiv*, abs/2211.09721, 2022.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and Ganguli, S. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pp. 2256–2265. PMLR, 2015.
- Song, J., Meng, C., and Ermon, S. Denoising diffusion implicit models. In *ICLR*, 2021a.
- Song, Y. and Ermon, S. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32, 2019.

- Song, Y. and Ermon, S. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020.
- Song, Y., Garg, S., Shi, J., and Ermon, S. Sliced score matching: A scalable approach to density and score estimation. In *Uncertainty in Artificial Intelligence*, pp. 574–584. PMLR, 2020.
- Song, Y., Sohl-Dickstein, J. N., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. *ArXiv*, abs/2011.13456, 2021b.
- Sun, J., Nie, W., Yu, Z., Mao, Z. M., and Xiao, C. Pointdp: Diffusion-driven purification against adversarial attacks on 3d point cloud recognition. *ArXiv*, abs/2208.09801, 2022.
- Tran, T., Do, T.-T., Reid, I. D., and Carneiro, G. Bayesian generative active deep learning. *ArXiv*, abs/1904.11643, 2019.
- Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., and Madry, A. Robustness may be at odds with accuracy. *arXiv: Machine Learning*, 2018.
- van den Oord, A., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding. *ArXiv*, abs/1807.03748, 2018.
- Wang, L., Xie, S., Li, T., Fonseca, R., and Tian, Y. Sample-efficient neural architecture search by learning action space. *ArXiv*, abs/1906.06832, 2019.
- Watson, D., Chan, W., Ho, J., and Norouzi, M. Learning fast samplers for diffusion models by differentiating through sample quality. In *ICLR*, 2022.
- Wei, C., Shen, K., Chen, Y., and Ma, T. Theoretical analysis of self-training with deep networks on unlabeled data. In *International Conference on Learning Representations*, 2020.
- Yang, L., Zhang, Z., and Hong, S. Diffusion models: A comprehensive survey of methods and applications. *arXiv* preprint arXiv:2209.00796, 2022.
- Zagoruyko, S. and Komodakis, N. Wide residual networks. *ArXiv*, abs/1605.07146, 2016.
- Zhang, C., Zhang, K., Zhang, C., Niu, A., Feng, J., Yoo, C. D., and Kweon, I.-S. Decoupled adversarial contrastive learning for self-supervised adversarial robustness. *ArXiv*, abs/2207.10899, 2022.
- Zhang, H. R., Yu, Y., Jiao, J., Xing, E. P., Ghaoui, L. E., and Jordan, M. I. Theoretically principled trade-off between robustness and accuracy. *ArXiv*, abs/1901.08573, 2019.

#### A. Theoretical Details for Section 3

#### A.1. Error probabilities in closed form.

Here, we briefly recapitulate the closed-form formulation for the standard and robust error probabilities as detailed in (Carmon et al., 2019; Deng et al., 2021).

The standard error probability can be written as

$$\operatorname{err}_{\operatorname{standard}}\left(f_{\boldsymbol{\theta}}\right) = \mathbb{P}\left(y \cdot \boldsymbol{x}^{\top} \boldsymbol{\theta} < 0\right) = \mathbb{P}\left(\mathcal{N}\left(\frac{\boldsymbol{\mu}^{\top} \boldsymbol{\theta}}{\sigma \|\boldsymbol{\theta}\|}, 1\right) < 0\right) = Q\left(\frac{\boldsymbol{\mu}^{\top} \boldsymbol{\theta}}{\sigma \|\boldsymbol{\theta}\|}\right), \tag{6}$$

where

$$Q(x) = \frac{1}{\sqrt{2\pi}} \int_{-\tau}^{\infty} e^{-t^2/2} dt$$

is the Gaussian error function and is non-increasing. Clearly the standard error probability is minimized when  $\frac{\theta}{\|\theta\|} = \frac{\mu}{\|\mu\|}$ , i.e.,  $\theta = c\mu$  for some scalar c > 0. We may impost  $\|\theta\|_2 = 1$  to ensure the unique solution  $\theta = \mu/\|\mu\|$ .

The robust error probability under the  $\ell_{\infty}$  adversarial set  $\Delta = \{\delta : \|\boldsymbol{\delta}\|_{\infty} \leq \epsilon\}$  is

$$\operatorname{err}_{\text{robust}}^{\infty,\varepsilon}(f_{\boldsymbol{\theta}}) = \mathbb{P}\left(\inf_{\|\boldsymbol{\nu}\|_{\infty} \leq \varepsilon} \left\{ y \cdot (\boldsymbol{x} + \boldsymbol{\nu})^{\top} \boldsymbol{\theta} \right\} < 0\right)$$

$$= \mathbb{P}\left( y \cdot \boldsymbol{x}^{\top} \boldsymbol{\theta} - \varepsilon \|\boldsymbol{\theta}\|_{1} < 0 \right) = \mathbb{P}\left( \mathcal{N}\left(\boldsymbol{\mu}^{\top} \boldsymbol{\theta}, \sigma^{2} \|\boldsymbol{\theta}\|^{2}\right) < \varepsilon \|\boldsymbol{\theta}\|_{1} \right)$$

$$= Q\left(\frac{\boldsymbol{\mu}^{\top} \boldsymbol{\theta}}{\sigma \|\boldsymbol{\theta}\|} - \frac{\varepsilon \|\boldsymbol{\theta}\|_{1}}{\sigma \|\boldsymbol{\theta}\|} \right). \tag{7}$$

In the following part, we use a simpler notation  $\operatorname{err}_{\operatorname{robust}}(f_{\theta})$  for the robust error  $\operatorname{err}_{\operatorname{robust}}^{\infty,\varepsilon}(f_{\theta})$  without ambiguity. The closed-form of the optimal  $\theta^*$  that minimizes the above robust error  $\operatorname{err}_{\operatorname{robust}}$  can be shown to be (Deng et al., 2021):

$$\boldsymbol{\theta}^* = \frac{T_{\varepsilon}(\boldsymbol{\mu})}{\|T_{\varepsilon}(\boldsymbol{\mu})\|},$$

where  $T_{\varepsilon}(\boldsymbol{\mu}): \mathbb{R}^d \to \mathbb{R}^d$  is the hard-thresholding operator with  $(T_{\varepsilon}(\boldsymbol{\mu}))_j = \mathrm{sign}\left(\boldsymbol{\mu}_j\right) \cdot \max\left\{\left|\boldsymbol{\mu}_j\right| - \varepsilon, 0\right\}$ . Under the mild assumption  $\boldsymbol{\mu}_j > \varepsilon, \forall j \in \{1, 2, \dots, d\}$ , the optimal solution can be simplified as:

$$oldsymbol{ heta}^* = rac{oldsymbol{\mu} - arepsilon \mathbf{1}_d}{\|oldsymbol{\mu} - arepsilon \mathbf{1}_d\|}.$$

Remark A.1. Note that when  $\mu = c\mathbf{1}_d$  for some constant  $c > \epsilon$ , the optimal solution  $\theta^* = \frac{\mu - \varepsilon \mathbf{1}_d}{\|\mu - \varepsilon \mathbf{1}_d\|}$  for minimizing the robust error is the same as the optimal solution  $\frac{\mu}{\|\mu\|}$  for minimizing the standard error. Otherwise, these two solutions are different, representing a trade-off between robustness and accuracy.

#### A.2. Details for the theoretical analysis in Section 3

Overall, we would like to design an appropriate synthetic distribution  $\widetilde{\mathcal{D}}$  that can help optimize the adversarial classification accuracy in the downstream task. First note that by Bayes rule, the optimal decision boundary for the true distribution  $\boldsymbol{x}|y\sim\mathcal{N}(y\boldsymbol{\mu},\sigma^2\mathbb{I})$  is given by  $\boldsymbol{\mu}^{\top}\boldsymbol{x}=0$ , i.e., the optimal classifier is parameterized by  $\boldsymbol{\theta}=c\boldsymbol{\mu}$  for any c>0. Therefore, we restrict our attention to synthetic data distributions that satisfy the following two conditions:

- 1. The marginal probability density  $p(\tilde{y})$  of the synthetic distribution matches p(y) of the real data distribution well.
- 2. The conditional probability densities  $p(\tilde{\boldsymbol{x}}|\tilde{y}=1)$  and  $p(\tilde{\boldsymbol{x}}|\tilde{y}=-1)$  of the synthetic data distribution are symmetric around the true optimal decision boundary  $\boldsymbol{\mu}^{\top}\boldsymbol{x}=0$ .

More specifically, we consider a special case of the synthetic data distribution  $\widetilde{\mathcal{D}}_{\mathcal{X}} = 0.5 \mathcal{N}(\tilde{\boldsymbol{\mu}}, \sigma^2 \mathbb{I}) + 0.5 \mathcal{N}(-\tilde{\boldsymbol{\mu}}, \sigma^2 \mathbb{I}).$ 

Proof of Proposition 3.2. We follow the proof strategy in (Carmon et al., 2019). Let  $b_i$  be the indicator that the i-th pseudo-label  $\tilde{y}_i$  assigned to  $\tilde{x}_i$  is incorrect, so that we have  $\tilde{x}_i \sim \mathcal{N}\left((1-2b_i)\,\tilde{y}_i\tilde{\mu},\sigma^2\mathbb{I}\right)$ . Let  $\gamma:=\frac{1}{\tilde{n}}\sum_{i=1}^{\tilde{n}}\left(1-2b_i\right)\in[-1,1]$  and  $\alpha:=\frac{\tilde{n}}{\tilde{n}+n}$ . Note that the true data samples  $x_i\sim\mathcal{N}\left(y_i\mu,\sigma^2\mathbb{I}\right)$ , thus we may write the final estimator as

$$\begin{split} \hat{\boldsymbol{\theta}}_{\text{final}} &= \frac{1}{n+\tilde{n}} (\sum_{j=1}^{\tilde{n}} \tilde{y}_{j} \tilde{\boldsymbol{x}}_{j} + \sum_{i=1}^{n} y_{i} \boldsymbol{x}_{i}) \\ &= \frac{1}{n+\tilde{n}} (\sum_{j=1}^{\tilde{n}} [(1-2b_{j}) \, \tilde{y}_{j} \tilde{\boldsymbol{\mu}} + \tilde{\boldsymbol{\epsilon}_{j}}] \cdot \tilde{y}_{j} + \sum_{i=1}^{n} [\boldsymbol{\mu} y_{i} + \boldsymbol{\epsilon}_{i}] \cdot y_{i}) \\ &= \alpha \gamma \tilde{\boldsymbol{\mu}} + \frac{1}{n+\tilde{n}} \sum_{i=1}^{\tilde{n}} \tilde{y}_{i} \tilde{\boldsymbol{\epsilon}}_{i} + (1-\alpha)\boldsymbol{\mu} + \frac{1}{n+\tilde{n}} \sum_{i=1}^{n} y_{i} \boldsymbol{\epsilon}_{i} \\ &= \alpha \gamma \tilde{\boldsymbol{\mu}} + (1-\alpha)\boldsymbol{\mu} + \frac{1}{n+\tilde{n}} (\sum_{i=1}^{n} y_{i} \boldsymbol{\epsilon}_{i} + \sum_{i=1}^{\tilde{n}} \tilde{y}_{i} \tilde{\boldsymbol{\epsilon}}_{i}), \end{split}$$

where  $\varepsilon_i$ ,  $\tilde{\varepsilon}_i \sim \mathcal{N}\left(\mathbf{0}, \sigma^2 \mathbb{I}\right)$  independent of each other, and the marginal probability density  $p(\tilde{y})$  matches p(y) well. Defining  $\tilde{\boldsymbol{\delta}} := \hat{\boldsymbol{\theta}}_{\text{final}} - \alpha \gamma \tilde{\boldsymbol{\mu}} - (1 - \alpha) \boldsymbol{\mu}$ .

By (6), we have that the standard error of  $f_{\hat{\theta}_{\text{final}}}$  is a non-increasing function of  $\frac{\mu^{\top}\hat{\theta}_{\text{final}}}{\sigma\|\hat{\theta}_{\text{final}}\|}$ . Note that when  $\tilde{n}$  is large enough, we have  $\alpha \to 1$  and the direction of  $\hat{\theta}_{\text{final}}$  approach the direction of  $\tilde{\mu}$ . Therefore, the statement in Case 1 holds as a consequence, and similarly for the robust error according to (7).

The remaining proof on Case 2 and Case 3 is based on a detailed discussion for the squared inverse of the term  $\frac{\mu^{\top} \hat{\theta}_{\text{final}}}{\sigma \|\hat{\theta}_{\text{final}}\|}$ :

$$\frac{\|\hat{\boldsymbol{\theta}}_{\text{final}}\|^2}{(\boldsymbol{\mu}^{\top}\hat{\boldsymbol{\theta}}_{\text{final}})^2} = \frac{\|\tilde{\boldsymbol{\delta}} + \alpha\gamma\tilde{\boldsymbol{\mu}} + (1-\alpha)\boldsymbol{\mu}\|^2}{(\alpha\gamma\langle\boldsymbol{\mu},\tilde{\boldsymbol{\mu}}\rangle + \boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}} + (1-\alpha)\|\boldsymbol{\mu}\|^2)^2}.$$
 (8)

Note that the larger the quantity in (8) is, the larger the standard error of  $f_{\hat{\theta}_{\text{final}}}$ .

Case 2. Assume  $\tilde{\mu} = c\mu$ . Then we have (8) reduces to:

$$\frac{\|\hat{\boldsymbol{\theta}}_{\text{final}}\|^{2}}{(\boldsymbol{\mu}^{\top}\hat{\boldsymbol{\theta}}_{\text{final}})^{2}} = \frac{\|\tilde{\boldsymbol{\delta}} + (1 - \alpha + c\gamma\alpha)\boldsymbol{\mu}\|^{2}}{\left((1 - \alpha + c\gamma\alpha)\|\boldsymbol{\mu}\|^{2} + \boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}}\right)^{2}}$$

$$= \frac{1}{\|\boldsymbol{\mu}\|^{2}} + \frac{\|\tilde{\boldsymbol{\delta}} + (1 - \alpha + c\gamma\alpha)\boldsymbol{\mu}\|^{2} - \frac{1}{\|\boldsymbol{\mu}\|^{2}}\left((1 - \alpha + c\gamma\alpha)\|\boldsymbol{\mu}\|^{2} + \boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}}\right)^{2}}{\left((1 - \alpha + c\gamma\alpha)\|\boldsymbol{\mu}\|^{2} + \boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}}\right)^{2}}$$

$$= \frac{1}{\|\boldsymbol{\mu}\|^{2}} + \frac{\|\tilde{\boldsymbol{\delta}}\|^{2} - \frac{1}{\|\boldsymbol{\mu}\|^{2}}(\boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}})^{2}}{\left((1 - \alpha + c\gamma\alpha)\|\boldsymbol{\mu}\|^{2} + \boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}}\right)^{2}},$$
(10)

which demonstrates that the bigger the c is, the smaller the standard error  $\operatorname{err}_{\operatorname{standard}}(f_{\hat{\boldsymbol{\theta}}_{\operatorname{final}}})$  is, which verifies the second part of Case 2.

Case 3. Assume  $\tilde{\mu} = c(\mu - \varepsilon \mathbf{1}_d)$ . Similar to Case 2, we rewrite the term inside the robust error function (7) as:

$$\frac{\|\hat{\boldsymbol{\theta}}_{\text{final}}\|^{2}}{\left((\boldsymbol{\mu} - \varepsilon \mathbf{1}_{d})^{\top} \hat{\boldsymbol{\theta}}_{\text{final}}\right)^{2}} = \frac{\|\tilde{\boldsymbol{\delta}} + (1 - \alpha)\boldsymbol{\mu} + c\gamma\alpha(\boldsymbol{\mu} - \varepsilon \mathbf{1}_{d})\|^{2}}{\left(c\gamma\alpha\|\boldsymbol{\mu} - \varepsilon \mathbf{1}_{d}\|^{2} + (1 - \alpha)\boldsymbol{\mu}^{\top}(\boldsymbol{\mu} - \varepsilon \mathbf{1}_{d}) + (\boldsymbol{\mu} - \varepsilon \mathbf{1}_{d})^{\top} \tilde{\boldsymbol{\delta}}\right)^{2}} 
\approx \frac{1}{\|\boldsymbol{\mu} - \varepsilon \mathbf{1}_{d}\|^{2}} + \frac{\|\tilde{\boldsymbol{\delta}}\|^{2} - \frac{1}{\|\boldsymbol{\mu} - \varepsilon \mathbf{1}_{d}\|^{2}} \left((\boldsymbol{\mu} - \varepsilon \mathbf{1}_{d})^{\top} \tilde{\boldsymbol{\delta}}\right)^{2}}{\left(c\gamma\alpha\|\boldsymbol{\mu} - \varepsilon \mathbf{1}_{d}\|^{2} + (\boldsymbol{\mu} - \varepsilon \mathbf{1}_{d})^{\top} \tilde{\boldsymbol{\delta}}\right)^{2}}, \tag{11}$$

where the last approximation is due to  $\tilde{n}$  sufficiently large and thus  $\alpha \approx 1$ . The above equation demonstrates the larger the c is, the smaller the robust error  $\operatorname{err}_{\operatorname{robust}}(f_{\hat{\theta}_{\operatorname{final}}})$  is, which proves the second part of Case 3.

*Proof of Theorem 3.4.* We follow the proof strategy in (Carmon et al., 2019), with the main difference being twofold: (i) in our final estimate  $\hat{\theta}_{\text{final}}$  depends on both the real data and synthetic data; (ii) our synthetic data are generated from a different distribution from the true distribution.

Below we analyze the sample complexity to achieve desired robust accuracy. Recall that, under the mild assumption  $\mu_j > \varepsilon, \forall j \in \{1, 2, \dots, d\}$ , the closed form of the robust error of  $f_{\theta}$  is  $Q(\frac{(\mu - \varepsilon \mathbf{1}_d)^{\top} \theta}{\sigma \|\theta\|})$ . Since  $\hat{\theta}_{\text{final}} = (1 - \alpha + c\gamma\alpha)\mu + \tilde{\delta}$ , we have the term inside  $Q(\cdot)$  function is:

$$\frac{(\boldsymbol{\mu} - \varepsilon \mathbf{1}_d)^{\top} \hat{\boldsymbol{\theta}}_{\text{final}}}{\sigma \|\hat{\boldsymbol{\theta}}_{\text{final}}\|} = \frac{(1 - \alpha + c\gamma\alpha)\|\boldsymbol{\mu}\|^2 + \boldsymbol{\mu}^{\top} \tilde{\boldsymbol{\delta}}}{\sigma \|\boldsymbol{\theta}\|} - \frac{(\varepsilon \mathbf{1}_d)^{\top} \hat{\boldsymbol{\theta}}_{\text{final}}}{\sigma \|\hat{\boldsymbol{\theta}}_{\text{final}}\|}.$$
 (12)

We consider the parameter setting:

$$\epsilon < \frac{1}{2}, \ \sigma = (nd)^{1/4}, \ \|\boldsymbol{\mu}\|^2 = d.$$
 (13)

Under such a setting and under the regime that  $d/n \gg 1$ , we have the classifier  $f_{\mu}$  achieves almost optimal performance in both robust and standard accuracy. Thus in the following, we mainly focus on the problem of finding the minimum number of synthetic samples  $\tilde{n}$  needed in order to ensure the estimate  $\hat{\theta}_{\text{final}}$  is close (in direction) to  $\mu$ .

For the squared inverse of the first term in (12), we have

$$\frac{\|\hat{\boldsymbol{\theta}}_{\text{final}}\|^{2}}{\left(\boldsymbol{\mu}^{\top}\hat{\boldsymbol{\theta}}_{\text{final}}\right)^{2}} = \frac{\|\tilde{\boldsymbol{\delta}} + (1 - \alpha + c\gamma\alpha)\boldsymbol{\mu}\|^{2}}{\left((1 - \alpha + c\gamma\alpha)\|\boldsymbol{\mu}\|^{2} + \boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}}\right)}$$

$$= \frac{1}{\|\boldsymbol{\mu}\|^{2}} + \frac{\|\tilde{\boldsymbol{\delta}}\|^{2} - \frac{1}{\|\boldsymbol{\mu}\|^{2}}\left(\boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}}\right)^{2}}{\left((1 - \alpha + c\gamma\alpha)\|\boldsymbol{\mu}\|^{2} + \boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}}\right)^{2}}$$

$$\leq \frac{1}{\|\boldsymbol{\mu}\|^{2}} + \frac{\|\tilde{\boldsymbol{\delta}}\|^{2}}{\left((1 - \alpha + c\gamma\alpha)\|\boldsymbol{\mu}\|^{2} + \boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}}\right)^{2}}$$

$$= \frac{1}{\|\boldsymbol{\mu}\|^{2}} + \frac{1}{\|\boldsymbol{\mu}\|^{4}} \frac{\|\tilde{\boldsymbol{\delta}}\|^{2}}{\left((1 - \alpha + c\gamma\alpha) + \frac{1}{\|\boldsymbol{\mu}\|^{2}}\boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}}\right)^{2}}$$
(14)

Note that due to the dependence between  $\tilde{y}_i$  and  $\tilde{\boldsymbol{\varepsilon}}_i$ , the random variable  $\tilde{\boldsymbol{\delta}}$  is non-Gaussian. To obtain the concentration bounds for  $\|\tilde{\boldsymbol{\delta}}\|^2$  and  $\boldsymbol{\mu}^\top \tilde{\boldsymbol{\delta}}$ , we follow the approach used in (Carmon et al., 2019) as follows. Recall  $\hat{\boldsymbol{\theta}}_{inter} = \frac{1}{n} \sum_{i=1}^n y_i \boldsymbol{x}_i$ ,  $\tilde{y}_i = \text{sign}(\hat{\boldsymbol{\theta}}_{inter}^\top \tilde{\boldsymbol{x}}_i)$ , and  $\tilde{\boldsymbol{\delta}} = \frac{1}{n+\tilde{n}} (\sum_{i=1}^n y_i \boldsymbol{\varepsilon}_i + \sum_{i=1}^{\tilde{n}} \tilde{y}_i \tilde{\boldsymbol{\varepsilon}}_i)$ . Find a coordinate system such that the first coordinate is in the direction of  $\hat{\boldsymbol{\theta}}_{inter}$ , and let  $v^{(i)}$  denote the i th entry of vector v in this coordinate system. Then

$$\tilde{y}_i = \operatorname{sign}\left(\tilde{\boldsymbol{x}}_i^{(1)}\right) = \operatorname{sign}\left(\boldsymbol{\mu}^{(1)} + \tilde{\boldsymbol{\varepsilon}}_i^{(1)}\right).$$

Consequently,  $\tilde{\boldsymbol{\varepsilon}}_{i}^{(j)}$  is independent of  $\tilde{y}_{i}$  for all i and  $j \geq 2$ , so that  $\tilde{y}_{i}\tilde{\boldsymbol{\varepsilon}}_{i}^{(j)} \sim \mathcal{N}\left(0,\sigma^{2}\right)$  and  $\frac{1}{n+\tilde{n}}(\sum_{i=1}^{n}y_{i}\boldsymbol{\varepsilon}_{i}^{(j)} + \sum_{i=1}^{\tilde{n}}\tilde{y}_{i}\tilde{\boldsymbol{\varepsilon}}_{i}^{(j)}) \sim \mathcal{N}\left(0,\sigma^{2}/(n+\tilde{n})\right)$  and

$$\sum_{j=2}^{d} \left( \frac{1}{n+\tilde{n}} \left( \sum_{i=1}^{n} y_i \boldsymbol{\varepsilon}_i^{(j)} + \sum_{i=1}^{\tilde{n}} \tilde{y}_i \tilde{\boldsymbol{\varepsilon}}_i^{(j)} \right) \right)^2 \sim \frac{\sigma^2}{n+\tilde{n}} \chi_{d-1}^2.$$

By the Cauchy-Schwarz inequality, we have:

$$\left(\frac{1}{n+\tilde{n}}(\sum_{i=1}^{n}y_{i}\varepsilon_{i}^{(1)} + \sum_{i=1}^{\tilde{n}}\tilde{y}_{i}\tilde{\varepsilon}_{i}^{(1)})\right)^{2} \\
\leq \frac{1}{(n+\tilde{n})^{2}} \left\{ \left(\sum_{i=1}^{n}y_{i}^{2} + \sum_{i=1}^{\tilde{n}}\tilde{y}_{i}^{2}\right) \left(\sum_{i=1}^{n}\left[\varepsilon_{i}^{(1)}\right]^{2} + \sum_{i=1}^{\tilde{n}}\left[\tilde{\varepsilon}_{i}^{(1)}\right]^{2}\right) \right\} \\
= \frac{1}{n+\tilde{n}}(\sum_{i=1}^{n}\left[\varepsilon_{i}^{(1)}\right]^{2} + \sum_{i=1}^{\tilde{n}}\left[\tilde{\varepsilon}_{i}^{(1)}\right]^{2}) \sim \frac{\sigma^{2}}{n+\tilde{n}}\chi_{n+\tilde{n}}^{2}.$$

Since  $\|\tilde{\pmb{\delta}}\|^2 = \sum_{j=1}^d \left(\frac{1}{n+\tilde{n}}(\sum_{i=1}^n y_i \pmb{\varepsilon}_i^{(j)} + \sum_{i=1}^{\tilde{n}} \tilde{y}_i \tilde{\pmb{\varepsilon}}_i^{(j)})\right)^2$ , we have by the union bound

$$\mathbb{P}\left(\|\tilde{\boldsymbol{\delta}}\|^{2} \geq 2\frac{\sigma^{2}}{n+\tilde{n}}(d-1+n+\tilde{n})\right) \leq \mathbb{P}\left(\chi_{n+\tilde{n}}^{2} \geq 2(n+\tilde{n})\right) + \mathbb{P}\left(\chi_{d-1}^{2} \geq 2(d-1)\right) < e^{-(\tilde{n}+n)/8} + e^{-(d-1)/8}.$$

Similarly applying the Cauchy-Schwarz inequality to  $\mu^{\top}\tilde{\delta} = \frac{1}{\tilde{n}+n} \left( \sum_{i=1}^{n} y_i \mu^{\top} \varepsilon_i + \sum_{i=1}^{\tilde{n}} \tilde{y}_i \mu^{\top} \tilde{\varepsilon}_i \right)$ , we have

$$\begin{split} & \left(\boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}}\right)^{2} \leq \frac{1}{(n+\tilde{n})^{2}} \left\{ \left(\sum_{i=1}^{n} y_{i}^{2} + \sum_{i=1}^{\tilde{n}} \tilde{y}_{i}^{2}\right) \left(\sum_{i=1}^{n} \left(\boldsymbol{\mu}^{\top}\boldsymbol{\varepsilon}_{i}\right)^{2} + \sum_{i=1}^{\tilde{n}} \left(\boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\varepsilon}}_{i}\right)^{2}\right) \right\} \\ & = \frac{1}{n+\tilde{n}} (\sum_{i=1}^{n} \left(\boldsymbol{\mu}^{\top}\boldsymbol{\varepsilon}_{i}\right)^{2} + \sum_{i=1}^{\tilde{n}} \left(\boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\varepsilon}}_{i}\right)^{2}) \sim \frac{\sigma^{2} \|\boldsymbol{\mu}\|^{2}}{n+\tilde{n}} \chi_{n+\tilde{n}}^{2}. \end{split}$$

Therefore we have

$$\mathbb{P}\left(|\boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}}| \geq \sqrt{2}\sigma\|\boldsymbol{\mu}\|\right) = \mathbb{P}\left((\boldsymbol{\mu}^{\top}\tilde{\boldsymbol{\delta}})^{2} \geq 2\sigma^{2}\|\boldsymbol{\mu}\|^{2}\right) \leq e^{-(\tilde{n}+n)/8}.$$

Finally, we look at the random variable  $\gamma$ . By definition  $\gamma = \frac{1}{\tilde{n}} \sum_{i=1}^{\tilde{n}} (1-2b_i)$ , where  $b_i$  is the indicator that  $\tilde{y}_i$  is incorrect for the feature  $\tilde{\boldsymbol{x}}_i$ . Denote  $\tilde{y}_i^{\circ} \in \{-1,1\}$  as the true label for  $\tilde{\boldsymbol{x}}_i$ , thus we have  $\tilde{\boldsymbol{x}}_i \sim N(c\tilde{y}_i^{\circ}\boldsymbol{\mu},\sigma^2)$ . Therefore

$$\mathbb{E}[b_i] = \mathbb{P}[b_i = 1] = \mathbb{P}\left(\tilde{y}_i^{\circ} \cdot \tilde{\boldsymbol{x}}_i^{\top} \hat{\boldsymbol{\theta}}_{\text{inter}} < 0\right) = \mathbb{P}\left(\mathcal{N}\left(\frac{c\boldsymbol{\mu}^{\top} \hat{\boldsymbol{\theta}}_{\text{inter}}}{\sigma \|\hat{\boldsymbol{\theta}}_{\text{inter}}\|}, 1\right) < 0\right) = Q\left(\frac{c\boldsymbol{\mu}^{\top} \hat{\boldsymbol{\theta}}_{\text{inter}}}{\sigma \|\hat{\boldsymbol{\theta}}_{\text{inter}}\|}\right).$$

Moreover since  $b_i$  are Bernoulli random variables, we have  $Var(b_i) = \mathbb{E}[b_i](1 - \mathbb{E}[b_i]) \leq \mathbb{E}[b_i]$ .

By definition of  $Q(\cdot)$  we clearly have  $Q\left(\frac{c\boldsymbol{\mu}^{\top}\hat{\boldsymbol{\theta}}_{\text{inter}}}{\sigma\|\hat{\boldsymbol{\theta}}_{\text{inter}}\|}\right) \leq Q\left(\frac{\boldsymbol{\mu}^{\top}\hat{\boldsymbol{\theta}}_{\text{inter}}}{\sigma\|\hat{\boldsymbol{\theta}}_{\text{inter}}\|}\right)$  when  $c \geq 1$  and  $\boldsymbol{\mu}^{\top}\hat{\boldsymbol{\theta}}_{\text{inter}} > 0$  (which happens with high probability as shown below). Thus

$$\mathbb{E}\left[\gamma\mid\hat{\boldsymbol{\theta}}_{\mathrm{inter}}\right] = 1 - 2Q\left(\frac{c\boldsymbol{\mu}^{\top}\hat{\boldsymbol{\theta}}_{\mathrm{inter}}}{\sigma\|\hat{\boldsymbol{\theta}}_{\mathrm{inter}}\|}\right) \geq 1 - 2\operatorname{err}_{\mathrm{standard}}\left(f_{\hat{\boldsymbol{\theta}}_{\mathrm{inter}}}\right),$$

where  $err_{standard}$  is given in (6).

Therefore, we expect  $\gamma$  to be reasonably large as long as  $\operatorname{err}_{\operatorname{standard}}(f_{\hat{\theta}_{\operatorname{inter}}}) < \frac{1}{2}$ . Similar to (Carmon et al., 2019), we have

$$\begin{split} \mathbb{P}\left(\gamma < \frac{1}{6}\right) &= \mathbb{P}\left(\frac{1}{\tilde{n}}\sum_{i=1}^{\tilde{n}}\left(1 - 2b_i\right) < \frac{1}{6}\right) \\ &= \mathbb{P}\left(\frac{1}{\tilde{n}}\sum_{i=1}^{\tilde{n}}\left(1 - 2b_i\right) < \frac{1}{6}|\operatorname{err}_{\operatorname{standard}}\left(f_{\hat{\boldsymbol{\theta}}_{\operatorname{inter}}}\right) > \frac{1}{3}\right) \cdot \mathbb{P}\left(\operatorname{err}_{\operatorname{standard}}\left(f_{\hat{\boldsymbol{\theta}}_{\operatorname{inter}}}\right) > \frac{1}{3}\right) \\ &+ \mathbb{P}\left(\frac{1}{\tilde{n}}\sum_{i=1}^{\tilde{n}}\left(1 - 2b_i\right) < \frac{1}{6}|\operatorname{err}_{\operatorname{standard}}\left(f_{\hat{\boldsymbol{\theta}}_{\operatorname{inter}}}\right) \leq \frac{1}{3}\right) \cdot \mathbb{P}\left(\operatorname{err}_{\operatorname{standard}}\left(f_{\hat{\boldsymbol{\theta}}_{\operatorname{inter}}}\right) \leq \frac{1}{3}\right) \\ &\leq \mathbb{P}\left(\operatorname{err}_{\operatorname{standard}}\left(f_{\hat{\boldsymbol{\theta}}_{\operatorname{inter}}}\right) > \frac{1}{3}\right) + \mathbb{P}\left(\frac{1}{\tilde{n}}\sum_{i=1}^{\tilde{n}}b_i > \frac{5}{12}|\operatorname{err}_{\operatorname{standard}}\left(f_{\hat{\boldsymbol{\theta}}_{\operatorname{inter}}}\right) \leq \frac{1}{3}\right). \end{split}$$

For the first probability, note that

$$\frac{1}{3} \ge Q\left(\frac{1}{2}\right) \ge Q\left(\left[2\left(1+\sqrt{n/d}\right)\right]^{-1/2}\right)$$

Therefore, by Lemma 1 in (Carmon et al., 2019), for sufficiently large d/n,

$$\mathbb{P}\left(\mathrm{err}_{\mathrm{standard}}\ \left(f_{\hat{\pmb{\theta}}_{\mathrm{inter}}}\right) > \frac{1}{3}\right) \leq e^{-c_1 \cdot \min\left\{\sqrt{d/n}, n(d/n)^{1/4}\right\}}$$

for some constant  $c_1$ .

For the second probability, note that  $b_i$  are i.i.d. Bernoulli random variables with mean value  $Q\left(\frac{c\mu^{\top}\hat{\theta}_{\text{inter}}}{\sigma\|\hat{\theta}_{\text{inter}}\|}\right) \leq \exp_{\text{standard}}\left(f_{\hat{\theta}_{\text{inter}}}\right)$ . Therefore, by Hoeffding's inequality we have

$$\mathbb{P}\left(\frac{1}{\tilde{n}}\sum_{i=1}^{\tilde{n}}b_i > \frac{5}{12}\mid \operatorname{err}_{\operatorname{standard}}\left(f_{\hat{\boldsymbol{\theta}}_{\operatorname{inter}}}\right) \leq \frac{1}{3}\right) \leq e^{-2\tilde{n}\left(\frac{5}{12} - \frac{1}{3}\right)^2} = e^{-\tilde{n}/72}.$$

Define the event,

$$\mathcal{E} := \left\{ \|\tilde{\boldsymbol{\delta}}\|^2 \leq 2 \frac{\sigma^2}{n + \tilde{n}} (d + n + \tilde{n}), \ \left| \boldsymbol{\mu}^\top \tilde{\boldsymbol{\delta}} \right| \leq \sqrt{2} \sigma \|\boldsymbol{\mu}\|, \ \text{and} \ \gamma \geq \frac{1}{6} \right\},$$

thus by the previous concentration bounds, we have

$$\mathbb{P}[\mathcal{E}^C] \le e^{-(\tilde{n}+n)/8} + e^{-(d-1)/8} + e^{-c_1 \cdot \min\left\{\sqrt{d/n}, n(d/n)^{1/4}\right\}} + e^{-\tilde{n}/72} \le e^{-c_2 \cdot \min\left\{\tilde{n}, \sqrt{d/n}, n(d/n)^{1/4}\right\}}.$$

Suppose the event  $\mathcal{E}$  holds, then for the formula in (14) we have:

$$\frac{\left\|\hat{\boldsymbol{\theta}}_{\text{final}}\right\|^2}{\left(\boldsymbol{\mu}^{\top}\hat{\boldsymbol{\theta}}_{\text{final}}\right)^2} \leq \frac{1}{\|\boldsymbol{\mu}\|^2} + \frac{2\sigma^2(d+n+\tilde{n})}{(n+\tilde{n})\|\boldsymbol{\mu}\|^4\left((1-\alpha+\frac{1}{6}c\alpha)-\frac{\sqrt{2}\sigma}{\|\boldsymbol{\mu}\|}\right)^2},$$

which, after substituting the parameter setting (13), translates into:

$$\frac{\sigma^{2} \left\| \hat{\boldsymbol{\theta}}_{\text{final}} \right\|^{2}}{\left( \boldsymbol{\mu}^{\top} \hat{\boldsymbol{\theta}}_{\text{final}} \right)^{2}} \leq \sqrt{\frac{n}{d}} + \frac{2nd(d+n+\tilde{n})}{(n+\tilde{n})d^{2} \left( (1-\alpha + \frac{1}{6}c\alpha) - \sqrt{2} \left( \frac{n}{d} \right)^{1/4} \right)^{2}} \\
\leq \sqrt{\frac{n}{d}} + \frac{2nd(d+n+\tilde{n})}{(n+\tilde{n})d^{2} \left( \frac{1}{6}c\alpha - \sqrt{2} \left( \frac{n}{d} \right)^{1/4} \right)^{2}} \\
\leq \sqrt{\frac{n}{d}} + \frac{72n}{c\tilde{n}} \left( 1 + \tilde{c}_{1} \left( \frac{n}{d} \right)^{1/4} \right),$$

where  $\tilde{c}_1$  is some positive constant, and above we also implicitly assumed that d/n is sufficiently large.

Combining the above results together, following the analysis in (Carmon et al., 2019), we conclude that there exists a universal constant  $\tilde{C}$  such that for  $\epsilon^2 \sqrt{d/n} \geq \tilde{C}$ , where n is the number of labeled real data used to construct the intermediate classifier, and additional  $\tilde{n}$  synthetic feature generated with mean vector  $\pm c\mu$  and pseudo labels, we have if

$$\tilde{n} \ge \frac{288n}{c} \epsilon^2 \sqrt{\frac{d}{n}},$$

we have

$$\mathbb{E}_{\hat{\boldsymbol{\theta}}_{\text{final}}} \text{ err } _{\text{robust}}^{\infty,\epsilon} \left( f_{\hat{\boldsymbol{\theta}}_{\text{final}}} \right) \leq Q \left( [\sqrt{2} - 1] \epsilon \left( d/n \right)^{1/4} \right) + e^{-\epsilon^2 c_2 \sqrt{d/n}} \leq 10^{-3}.$$

for sufficiently large  $\tilde{C}$ .

# B. More simulation results under Gaussian setting in Section 3

In this section, we present more detailed simulation results under the Gaussian setting in Section 3 to demonstrate different scenarios in Proposition 3.2.

**Experimental setting** The experimental pipeline is as follow: 1) learning a intermediate classifier  $\hat{\theta}_{\text{inter}}$  by n label data  $\mathcal{D} = \{(\boldsymbol{x}_1, y_1), \dots, (\boldsymbol{x}_n, y_n)\}$ , 2) generating  $\tilde{n}$  synthetic data  $\tilde{\boldsymbol{x}} \sim \widetilde{\mathcal{D}}_{\mathcal{X}} = 0.5 \mathcal{N}(\tilde{\boldsymbol{\mu}}, \sigma^2 \mathbb{I}) + 0.5 \mathcal{N}(-\tilde{\boldsymbol{\mu}}, \sigma^2 \mathbb{I})$ , with  $\tilde{\boldsymbol{\mu}} = c\boldsymbol{\mu}$ , 3) assigning pseudo label for synthetic data using the intermediate classifier  $\hat{\theta}_{\text{inter}}$ , 4) learning  $\hat{\theta}_{\text{final}}$  by adversarial training on  $\tilde{n}$  synthetic data, 5) testing on 10k extra real data to obtain the clean accuracy and robust accuracy. For data dimension d = 100, we set  $\|\boldsymbol{\mu}\|^2 = 2$ ,  $\varepsilon = 0.5$ , and for d = 100, we set  $\|\boldsymbol{\mu}\|^2 = 4$ ,  $\varepsilon = 0.1$ .

Table 4 and Table 5 show the clean and robust accuracy learned on synthetic distribution  $\tilde{\mu} = c\mu$  with different angles between  $\mu$  and  $\epsilon 1_d$ . Table 7 shows the clean and robust accuracy learned on synthetic distribution  $\tilde{\mu} = c(\mu - \varepsilon 1_d)$  with different angles between  $\mu$  and  $\epsilon 1_d$ . Recall that  $\mu$  is (one of) the optimal linear classifier that maximizes the clean accuracy under the true distribution considered in Section 3, similarly  $\mu - \epsilon 1_d$  is the optimal solution for robust accuracy. Therefore, different angles between  $\mu$  and  $\epsilon 1_d$  represent different trade-offs between the clean and robust accuracy. For example, when the angle between  $\mu$  and  $\epsilon 1_d$  is 0 degrees, i.e.,  $\mu = c 1_d$ , we have that the optimal solution for clean accuracy and robust accuracy are the same. In most cases, the classifier learned from the synthetic distribution that is most separable achieves better performance even than the iid samples, which verifies Proposition 3.2.

Table 4. The clean and robust accuracy learned on synthetic distribution  $\tilde{\mu}=c\mu$  when d=2 and the angle between  $\mu$  and  $\epsilon$  is 0 degrees and 90 degrees. "Real" denotes the real data distribution, and n denotes the number of data from the real distribution, while we use "c" to denote different synthetic distributions and use  $\tilde{n}$  to denote the number of synthetic data. The results and the standard deviation in the bracket are obtained from 50 repetitions.

		0 de	gree	90 d	egree
		acc (std)	rob acc (std)	acc (std)	rob acc (std)
	n = 10	0.9201 (0.0012)	<b>0.7593</b> (0.0020)	0.9171 (0.0046)	0.7552 (0.0040)
Real	n = 20	0.9204 (0.0007)	<b>0.7598</b> (0.0016)	0.9186 (0.0017)	0.7563 (0.0012)
Real	n = 50	0.9206 (0.0004)	<b>0.7605</b> (0.0007)	0.9196 (0.0009)	0.7566 (0.0006)
	n = 100	0.9205 (0.0004)	<b>0.7608</b> (0.0006)	0.9199 (0.0006)	0.7565 (0.0007)
	$\tilde{n} = 10$	0.9159 (0.0099)	0.7541 (0.0096)	0.9104 (0.0121)	0.7492 (0.0122)
- 05	$\tilde{n} = 20$	0.9179 (0.0047)	0.7562 (0.0050)	0.9161 (0.0052)	0.7546 (0.0054)
c = 0.5	$\tilde{n} = 50$	0.9200 (0.0023)	0.7586 (0.0024)	0.9183 (0.0022)	0.7570 (0.0022)
	$\tilde{n} = 100$	0.9213 (0.0011)	0.7601 (0.0009)	0.9193 (0.0012)	0.7576 (0.0010)
	$\tilde{n} = 10$	0.9133 (0.0066)	0.7502 (0.0061)	0.9161 (0.0048)	<b>0.7598</b> (0.0048)
c = 1	$\tilde{n}=20$	0.9155 (0.0020)	0.7516 (0.0019)	0.9180 (0.0017)	<b>0.7612</b> (0.0020)
c = 1	$\tilde{n} = 50$	0.9161 (0.0009)	0.7525 (0.0006)	0.9186 (0.0010)	<b>0.7620</b> (0.0006)
	$\tilde{n} = 100$	0.9165 (0.0005)	0.7528 (0.0006)	0.9189 (0.0005)	<b>0.7622</b> (0.0003)
	$\tilde{n} = 10$	<b>0.9209</b> (0.0038)	0.7523 (0.0025)	<b>0.9221</b> (0.0017)	0.7583 (0.0015)
c = 1.5	$\tilde{n} = 20$	<b>0.9228</b> (0.0010)	0.7536 (0.0006)	<b>0.9226</b> (0.0013)	0.7588 (0.0013)
c = 1.0	$\tilde{n} = 50$	<b>0.9229</b> (0.0008)	0.7538 (0.0005)	<b>0.9232</b> (0.0005)	0.7594 (0.0006)
	$\tilde{n} = 100$	<b>0.9232</b> (0.0003)	0.7538 (0.0005)	<b>0.9233</b> (0.0005)	0.7595 (0.0005)

Table 5. The clean and robust accuracy learned on synthetic distribution  $\tilde{\mu}=c\mu$  when d=2 and the angle between  $\mu$  and  $\epsilon$  is 30 degrees and 60 degrees. "Real" denotes the real data distribution, and n denotes the number of data from the real distribution, while we use "c" to denote different synthetic distributions and use  $\tilde{n}$  to denote the number of synthetic data. The results and the standard deviation in the bracket are obtained from 50 repetitions.

		30 de	egree	60 d	egree
		acc (std)	rob acc (std)	acc (std)	rob acc (std)
	n = 10	0.8307 (0.0123)	0.6343 (0.0283)	0.8348 (0.0117)	0.6378 (0.0293)
Real	n = 20	0.8353 (0.0055)	0.6404 (0.0234)	0.8391 (0.005)	0.6433 (0.0222)
Keai	n = 50	0.8371 (0.0022)	0.6450 (0.0168)	0.8410 (0.0017)	0.6494 (0.0134)
	n = 100	0.8385 (0.0010)	0.6461 (0.0097)	0.8413 (0.0013)	0.6522 (0.0102)
	$\tilde{n} = 10$	0.8265 (0.0184 )	0.6282 (0.0418 )	0.8338 (0.0132 )	0.6303 (0.0335 )
- 05	$\tilde{n}=20$	0.8299 (0.0129)	0.6352 (0.0325)	0.8365 (0.0132)	0.6393 (0.0316)
c = 0.5	$\tilde{n} = 50$	0.8372 (0.0046)	0.6483 (0.0215)	0.8414 (0.0034)	0.6489 (0.0199)
	$\tilde{n} = 100$	0.8402 (0.0015)	0.6466 (0.0110)	0.8431 (0.0012)	0.6510 (0.0135)
	$\tilde{n} = 10$	0.8383 (0.0158)	0.6439 (0.0319)	<b>0.8377</b> (0.0074)	0.6396 (0.0267)
c = 1	$\tilde{n}=20$	0.8425 (0.0060)	0.6480 (0.0218)	<b>0.8416</b> (0.0034)	<b>0.6513</b> (0.0178)
c = 1	$\tilde{n} = 50$	0.8455 (0.0023)	0.6553 (0.0128)	<b>0.8432</b> (0.0020)	<b>0.6503</b> (0.0122)
	$\tilde{n} = 100$	0.8457 (0.0021)	0.6535 (0.0100)	<b>0.8435</b> (0.0014)	<b>0.6501</b> 1 (0.0096)
	$\tilde{n} = 10$	<b>0.8431</b> (0.0045)	<b>0.6542</b> (0.0173)	0.8368 (0.0073)	<b>0.6446</b> (0.0213)
c = 1.5	$\tilde{n}=20$	<b>0.8447</b> (0.0021)	<b>0.6542</b> (0.0142)	0.8393 (0.0022)	0.6479 (0.0150)
c = 1.5	$\tilde{n} = 50$	<b>0.8455</b> (0.0006)	<b>0.6556</b> (0.0082)	0.8404 (0.0005)	0.6488 (0.0089)
	$\tilde{n} = 100$	<b>0.8457</b> (0.0004)	<b>0.6547</b> (0.0057)	0.8404 (0.0007)	0.6486 (0.0082)

Table 6. The clean and robust accuracy learned on synthetic distribution  $\tilde{\mu}=c\mu$  when d=100 and the angle between  $\mu$  and  $\epsilon$  is 0 degrees. "Real" denotes the real data distribution, and n denotes the number of data from the real distribution, while we use "c" to denote different synthetic distributions and use  $\tilde{n}$  to denote the number of synthetic data. The results and the standard deviation in the bracket are obtained from 50 repetitions

	acc (std)	rob acc (std)
Real	0.9023 (0.0192) 0.9341 (0.0128) 0.9599 (0.0028) 0.9682 (0.0014)	0.6843 (0.0359) 0.7519 (0.0267) 0.8078 (0.0061) 0.8239 (0.0028)
c = 0.5	0.7562 (0.0564) 0.8566 (0.0307) 0.9261 (0.0117) 0.9505 (0.0047)	0.4611 (0.0694) 0.6047 (0.0491) 0.7328 (0.0227) 0.7848 (0.0111)
c = 1	0.8866 (0.0273) 0.9371 (0.0091) 0.9620 (0.0028) 0.9695 (0.0012)	0.6557 (0.0487) 0.7555 (0.0201) 0.8085 (0.0060) 0.8239 (0.0031)
c = 1.5	<b>0.9400</b> (0.0100) <b>0.9591</b> (0.0037) <b>0.9710</b> (0.0013) <b>0.9743</b> (0.0008)	<b>0.7603</b> (0.0233) <b>0.8031</b> (0.0080) <b>0.8280</b> (0.0028) <b>0.8343</b> (0.0011)

Table 7. The clean and robust accuracy learned on synthetic distribution  $\tilde{\mu}=c(\mu-\varepsilon\mathbf{1}_d)$  when d=2 and the angle between  $\mu$  and  $\epsilon$  is 30 degrees and 60 degrees. "Real" denotes the real data distribution, and n denotes the number of data from the real distribution, while we use " $\epsilon$ " to denote different synthetic distributions and use  $\tilde{n}$  to denote the number of synthetic data. The results and the standard deviation in the bracket are obtained from 50 repetitions.

	30 de	egree	60 de	egree
	acc (std)	rob acc (std)	acc (std)	rob acc (std)
Real	0.9152 (0.0049) 0.9170 (0.0030) 0.9183 (0.0009) 0.9185 (0.0006)	0.7633 (0.0111) 0.7642 (0.0075) 0.7653 (0.0040) 0.7658 (0.0027)	<b>0.9211</b> (0.0034) <b>0.9225</b> (0.0020) <b>0.9232</b> (0.0011) <b>0.9235</b> (0.0009)	0.7702 (0.0094) 0.7714 (0.0068) 0.7711 (0.0050) 0.7724 (0.0027)
c = 0.5	0.9089 (0.0183) 0.9144 (0.0068) 0.9174 (0.0029) 0.9183 (0.0016)	0.7563 (0.0310) 0.7659 (0.0107) 0.7680 (0.0068) 0.7681 (0.0053)	0.9111 (0.0114) 0.9138 (0.0068) 0.9161 (0.0038) 0.9165 (0.0031)	0.7638 (0.0172) 0.7694 (0.0066) 0.7714 (0.0033) 0.7727 (0.0014)
c = 1	0.9135 (0.0116) 0.9178 (0.0046) 0.9183 (0.0042) 0.9196 (0.0017)	0.7642 (0.0194) 0.7710 (0.0073) 0.7728 (0.0042) 0.7733 (0.0036)	0.9069 (0.0111) 0.9042 (0.0098) 0.9073 (0.0047) 0.9059 (0.0039)	0.7677 (0.0109) 0.7676 (0.0072) 0.7702 (0.0017) 0.7698 (0.0016)
c = 1.5	<b>0.9181</b> (0.0079) <b>0.9209</b> (0.0053) <b>0.9218</b> (0.0029) <b>0.9222</b> (0.0017)	<b>0.7747</b> (0.0104) <b>0.7770</b> (0.0052) <b>0.7788</b> (0.0028) <b>0.7793</b> (0.0023)	0.9034 (0.0079) 0.9077 (0.0059) 0.9073 (0.0030) 0.9077 (0.0024)	<b>0.7704</b> (0.0053) <b>0.7716</b> (0.0056) <b>0.7722</b> (0.0014) <b>0.7729</b> (0.0011)

#### C. The detailed construction of the contrastive loss

In this section, we first give a detailed description of several possible ways to design contrastive loss, especially in constructing positive and negative pairs. Then, we give a visualization of the synthetic data distributions generated under different contrastive losses.

#### C.1. Positive and negative pair selection strategy.

In this subsection, we give several possible ways to construct positive and negative pairs.

- 1. Vanilla version: Using all the samples in the minibatch is the common strategy for contrastive learning. In the diffusion process, since for each time step t, we want to distinguish each image from other images in the minibatch at the same time step, a straight-forward strategy is to use all the samples in the minibatch other than  $x_t^i$  at time step t to be the negative pairs. For the positive pairs, we can simply adopt  $x_{t+1}^i$  to be the positive pairs rather than augmentation of  $x_t^i$ .
- 2. Real data as positive pairs: A possible improvement upon the vanilla version is considering we aim to generate images similar to real data. Therefore, we can directly adopt the real data as positive pairs.
- 3. Real data as negative pairs: Another improvement upon the vanilla version is considering the other images in time step t in the minibatch is not as high quality as the real data. Therefore, we can directly adopt the real data as the negative pairs.
- 4. Class conditional version: When we use conditional diffusion, and the class label of  $x_t$  in the minibatch is available, a further improvement can be adopted is to use all the samples with different class label y in the minibatch at time step t to be the negative pairs.

#### C.2. Visualization of the synthetic data distribution generated by different designs of the contrastive loss

In this subsection, we demonstrate the synthetic distributions generated by different designs of the contrastive loss mentioned in Section C.1 on the Gaussian setting mentioned in Section 3.1. Figure 6 shows the synthetic distribution generated by using  $\mathcal{N}(\mathbf{0}, \mathbb{I})$  as initialization, while Figure 7 shows the synthetic distribution generated by using  $\mathcal{N}(\mathbf{0}, 4\mathbb{I})$  as initialization. In all figures, all of the contrastive loss except for conditional hard negative mining form a circle within each class, which means these algorithms cannot explicitly distinguish the data within the same class and thus maximize the distance within each class, while the guidance from conditional hard negative mining can generate samples that are more distinguishable.

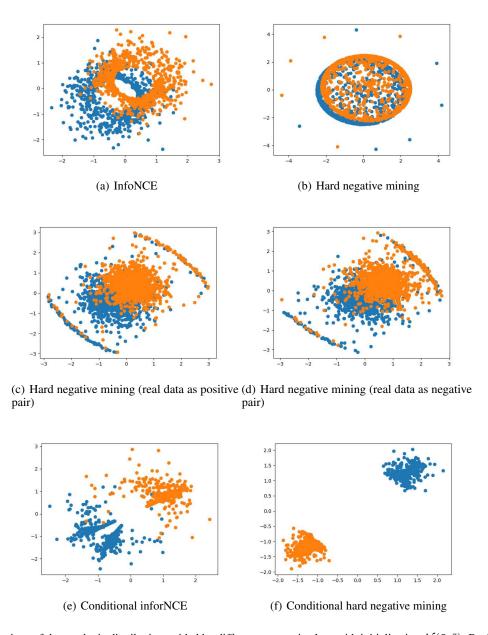


Figure 6. A comparison of the synthetic distribution guided by different contrastive loss with initialization  $\mathcal{N}(\mathbf{0}, \mathbb{I})$ . Real data as positive pair means using the mixture of oracle distribution  $\mathcal{N}(\pm \mathbf{1}_d, \mathbb{I})$  and the data in the same batch as negative pair, while real data as negative pair means using the data in the same batch as positive pair and using the mixture of oracle distribution as negative pair.

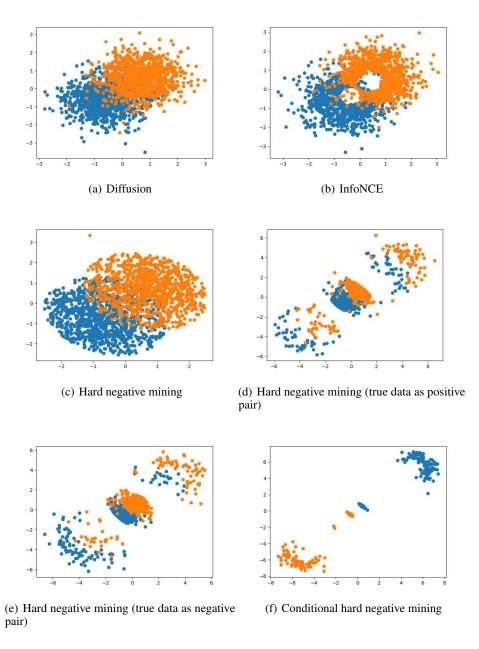


Figure 7. A comparison of the synthetic distribution guided by different contrastive loss with initialization  $\mathcal{N}(\mathbf{0}, 4\mathbb{I})$ .

# D. The experimental results for the real datasets

### D.1. Experimental setup for MNIST dataset

We describe the pipeline of synthetic data generation for adversarial robustness and a corresponding setting for the MNIST dataset in this subsection.

**Dataset.** MNIST dataset (LeCun et al., 1998) contains 60k 28×28 pixel grayscale handwritten digits between 0 to 9 for training and 10k digits for testing.

Synthetic data generation by the diffusion model. To utilize the pre-trained diffusion model  $^4$ , we use a conditional DDPM for generating samples for MNIST dataset. We adopt the hard negative mining loss with  $\tau=10$ , the strength of guidance of the contrastive loss  $\lambda=5k$ , the probability of the same class in the minibatch  $\tau^+=0.1$  and the hardness of negative mining  $\beta=1$ . We also use the pre-trained four layers Convolutional Neural Network model (removing the last fully connected layer) to get the representation for applying the contrastive loss and use a 2-layer feed-forward neural network to encode the representation after the pre-trained model.

**Adversarial Training.** Since grayscale handwritten digits can be easily classified, we adopt four layers Convolutional Neural Network as the classifier instead of using the Wide ResNet-28-10 model. We adopt stochastic weight averaging (Izmailov et al., 2018) with the decay rate 0.995 and use TRADES (Zhang et al., 2019) with 10 Projected Gradient Descent steps and  $\varepsilon_{\infty}=0.3$  for 150 epochs with batch size 1024.

# D.2. Experimental setup for CIFAR-10 dataset

We describe the pipeline of synthetic data generation for adversarial robustness and a corresponding setting for the CIFAR-10 dataset in this subsection.

**Dataset.** CIFAR-10 dataset (Krizhevsky, 2009) contains 50K 32x32 color training images in 10 classes and 10K images for testing.

**Overall training pipeline** We follow the same training pipeline as (Gowal et al., 2021), i.e., synthesizing data by using the diffusion model, assigning pseudo-label for synthetic data and aggregating the original data and the synthetic data for adversarial training. We give a careful explanation of these three components as follow.

Synthetic data generation by the diffusion model. Considering the advantage of DDIM on generation speed, we base on the official implementation of the DDIM model (Song et al., 2021a) and add the guidance of the contrastive loss. We generate images with 200 steps with batchsize=512, and use the quadratic version of sub-sequence selection  $^5$ . For the guidance of the contrastive loss, we try different designs of the contrastive loss mentioned in Section 4.2. We set the temperature  $\tau=0.1$  and the strength of guidance of the contrastive loss  $\lambda=20k$  in the InfoNCE loss, while  $\tau=10$ , the strength of guidance of the contrastive loss  $\lambda=100k$ , the probability of the same class in the minibatch  $\tau^+=0.1$  and the hardness of negative mining  $\beta=1$  in hard negative mining loss. These corresponding hyperparameters are chosen based on some preliminary experiments on image generation. The detailed ablation studies can be found in Section 5.2. Moreover, we also delve into the representation used by contrastive loss. The default setting is to use the pre-trained Wide ResNet-28-10 model (Gowal et al., 2021) to get the representation for applying the contrastive loss, which is named as (without embedding) in Section 5.2. A further improvement is to apply a 2-layer feed-forward neural network to encode the representation after the pre-trained model, which is named as (with embedding). The advantage of the latter design is we can adopt the contrastive loss to optimize the encoding network rather than a fixed encoder.

**LaNet for assigning pseudo-label.** Since the DDIM is an unconditional generator, we need to assign the pseudo-label to the generated sample. We follow the same choice adopted by (Sehwag et al., 2022), i.e., using state-of-the-art LaNet (Wang et al., 2019) network for assigning the pseudo-label for the synthetic data.

<sup>&</sup>lt;sup>4</sup>https://github.com/VSehwag/minimal-diffusion

<sup>&</sup>lt;sup>5</sup>We refer to Appendix D.2 for a detailed explanation of the quadratic version.

**Adversarial Training.** We follow the same setting as (Gowal et al., 2021), i.e., we use Wide ResNet-28-10 (Zagoruyko & Komodakis, 2016) with Swish activation function (Hendrycks & Gimpel, 2016), adopt stochastic weight averaging (Izmailov et al., 2018) with decay rate 0.995 and use TRADES (Zhang et al., 2019) with 10 Projected Gradient Descent steps and  $\varepsilon_{\infty} = 8/255$  for 400 epochs with batch size  $1024^6$ .

**Evaluation setup** For each trained model, we adopt AUTOATTACK (Croce & Hein, 2020) with  $\epsilon_{\infty} = 8/255$ .

# D.3. Experimental setup for Traffic Signs dataset

We describe the pipeline of synthetic data generation for adversarial robustness and a corresponding setting for the Traffic Signs dataset in this subsection.

**Dataset.** Traffic Signs dataset (Houben et al., 2013) contains 39252 training images in 43 classes and 12629 images for testing, and the image sizes vary between 15x15 to 250x250 pixels.

Synthetic data generation by the diffusion model. To utilize the pre-trained diffusion model  $^7$ , we use a conditional DDPM for generating samples for Traffic Signs dataset. We adopt the hard negative mining loss with  $\tau=10$ , the strength of guidance of the contrastive loss  $\lambda=5k$ , the probability of the same class in the minibatch  $\tau^+=0.1$  and the hardness of negative mining  $\beta=1$ . We also use the pre-trained Wide ResNet-28-10 model to get the representation for applying the contrastive loss and use a 2-layer feed-forward neural network to encode the representation after the pre-trained model.

**Adversarial Training.** We follow the same setting as the CIFAR-10 dataset, except the training epochs are reduced to 50. We also extend the training epochs to 400 but do not find significant improvement.

# E. Ablation study<sup>8</sup>.

#### E.1. The effectiveness of different contrastive losses.

Table 8 demonstrates the performance of different designs of the contrastive loss. We find out that applying the hard negative mining together with the embedding network achieves better clean and robust accuracy when the additional data is small (50K and 200K setting), while the infoNCE loss achieves better clean and robust accuracy when the additional data is large (1M setting). This result shows that we can improve the sample efficiency of the generative model by carefully designing the contrastive loss.

Table 8. The performance of Contrastive-DP under different contrastive loss: infoNCE and HNM losses, and w/wo embedding denote with/without an embedding network.

	50K		200K		1M	
	clean acc	rob acc	clean acc	rob acc	clean acc	rob acc
DDIM+infoNCE	83.40%	52.74%	84.18%	54.75%	85.64%	56.28%
DDIM+HNM(w embedding)	84.20%	53.19%	85.71%	54.92%	85.29%	56.12%
DDIM+HNM(wo embedding)	83.97%	52.89%	85.65%	54.83%	85.38%	55.95%

#### E.2. Sensitivity of the strength of the contrastive loss

Table 9 shows the influence of the strength of the contrastive loss.  $\lambda = 100k$  gives consistently better results than a smaller  $\lambda = 50k$  or a larger  $\lambda = 200k$  on robust accuracy on all settings. Moreover, we find the larger the  $\lambda$  is, the better performance we get on clean accuracy when the additional data is small (50K case), while the smaller the  $\lambda$  is, the better performance we get on clean accuracy when the additional data is large (1M case).

<sup>&</sup>lt;sup>6</sup>For Table 8 in the ablation studies subsection, we use batch size with 256.

<sup>&</sup>lt;sup>7</sup>https://github.com/VSehwag/minimal-diffusion

<sup>&</sup>lt;sup>8</sup>In this section, the robust accuracy is reported by the worst accuracy obtained by either AUTOATTACK (Croce & Hein, 2020) or AA+MT (Gowal et al., 2020)

*Table 9.* The performance of Contrastive-DP under different  $\lambda$  values.

	50K		200K		1M	
	clean acc	rob acc	clean acc	rob acc	clean acc	rob acc
$\lambda = 50k$	84.41%	53.78%	85.45%	55.24%	86.35%	56.83%
$\lambda = 100k$	83.66%	53.91%	85.71%	55.79%	86.30%	56.84%
$\lambda = 200k$	84.51%	53.55%	85.51%	55.33%	85.98%	56.69%

#### E.3. Data selection for synthetic data

Data selection methods are worthy of study since, in practice, we would like to know whether we can achieve better performance by generating a large number of samples and applying some selection criteria to filter out some samples. Therefore, we propose several data selection criterion and evaluate corresponding effectiveness in Table 10. All of the selection methods on Contrastive-DP are higher than vanilla DDIM plus selection methods, which demonstrates the superiority of using the contrastive learning loss as the guidance rather than using selection methods on the images generated by the vanilla diffusion model.

Table 10. Comparison of different data selection criteria. The detailed explanation of each selection method can be found in Append E.3.

	50K		200K		1M	
	clean acc	rob acc	clean acc	rob acc	clean acc	rob acc
DDIM (Separability)	79.93%	49.49%	85.09%	54.90%	84.87%	56.08%
Contrastive-DP (Gradient norm)	80.41%	49.47%	84.64%	55.17%	86.36%	57.11%
Contrastive-DP (Gradient norm-rob)	83.91%	55.23%	84.78%	55.42%	85.93%	57.18%
Contrastive-DP (Entropy)	83.66%	53.91%	85.71%	55.79%	86.30%	56.84%

Below we summarize different data selection methods:

- DDIM (Separability): We adopt the separability of the data as a criterion to make the selection of the data generated by vanilla DDIM. For each data, we use a pre-trained WRN-28-10 model to encode them into the embedding space. Then, we compute the L2 distance between each sample and the centroid of all classes (which is easily computed as the mean of all samples in this class) and add them together. To select a subset of samples that are most distinguishable, we choose the top K samples that have the smallest distance in each class.
- Contrastive-DP (Gradient norm): We use the gradient norm with respect to a pre-trained WRN-28-10 model as a criterion to make the selection on the data generated by Contrastive-DP. The larger the gradient norm is, the more informative the sample is for learning a downstream model. Therefore, we select the top K samples that have the largest gradient norm in each class.
- Contrastive-DP (Gradient norm-rob): Similar to Contrastive-DP (Gradient norm), we use the gradient norm of the robust loss rather than standard classification loss as a criterion to make the selection on the data generated by Contrastive-DP. Therefore, we select the top K samples that have the largest gradient norm in each class.
- Contrastive-DP (Entropy): We use the entropy of each sample with respect to LaNet as a criterion to make the selection on the data generated by Contrastive-DP. The smaller the entropy is, the higher likelihood this image has good quality. Therefore, we select the top K samples that have the smallest entropy in each class.

# F. Additional experiments

#### F.1. Changing the base adversarial training algorithm

We mainly adopt the TRADES (Zhang et al., 2019) for adversarial training on synthetic data together with real training data. A question is whether Contrastive-DP algorithm can also have good performance using vanilla adversarial training algorithm (Madry et al., 2017). Table 11 demonstrates Contrastive-DP also shows advantages against vanilla DDPM and DDIM by different base adversarial training algorithms.

Table 11. The clean and adversarial accuracy on CIFAR-10 dataset. The robust accuracy is reported by AUTOATTACK (Croce & Hein, 2020) with  $\epsilon_{\infty}=8/255$  and WRN-28-10. 50k, 200k, and 1M denote the number of synthetic used for adversarial training.

	50K		200	200K		1M	
	clean acc	rob acc	clean acc	rob acc	clean acc	rob acc	
DDIM Contrastive-DP	87.84% 88.50%	<b>54.97</b> % 54.74%	<b>89.19%</b> 88.26%	53.79% <b>54.20</b> %	88.91% <b>89.43</b> %	55.10% <b>55.31%</b>	
DDPM Contrastive-DP	88.19% <b>88.99</b> %	53.32% <b>53.67</b> %	89.21% <b>89.55</b> %	54.16% <b>54.85</b> %	<b>89.98%</b> 89.97%	54.17% <b>55.82%</b>	

#### F.2. Comparison with adversarial self-supervised learning

In the main paper, we only give the comparison of Contrastive-DP with the state-of-the-art method of adversarial robustness (Gowal et al., 2021) by using the diffusion model to generate synthetic data. Since contrastive learning is also used in adversarial self-supervised learning literature (Kim et al., 2020a; Fan et al., 2021; Zhang et al., 2022), we give a detailed comparison with these methods in Table 12, which also demonstrates the effectiveness of Contrastive-DP.

Table 12. The clean and adversarial accuracy on CIFAR-10 dataset. The robust accuracy is reported by AUTOATTACK (Croce & Hein, 2020) with  $\epsilon_{\infty}=8/255$ .

	rob acc
RoCL (Kim et al., 2020a)	47.88%
AdvCL (Fan et al., 2021)	49.77%
DeACL (Zhang et al., 2022)	50.39%
Contrastive-DP	59.99%