# Label-Noise Robust Generative Adversarial Networks

Takuhiro Kaneko[1]    Yoshitaka Ushiku[1]    Tatsuya Harada[1,2]

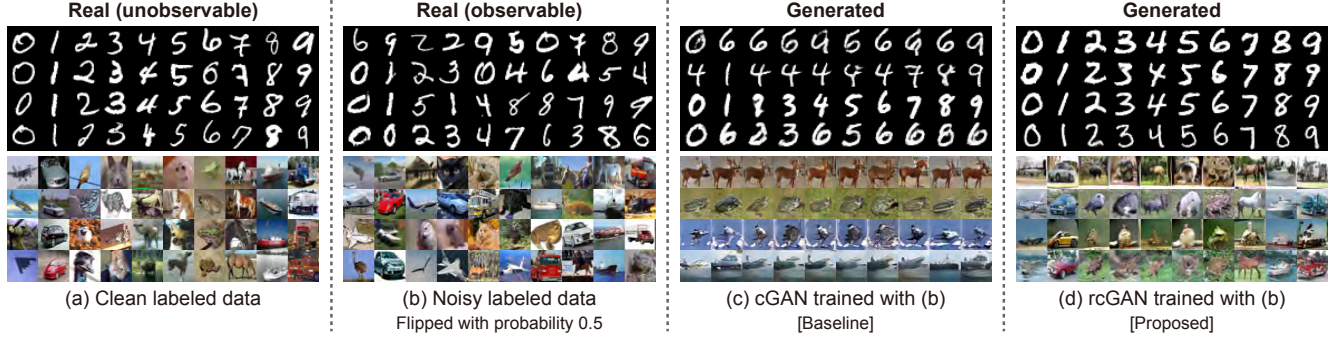[1]The University of Tokyo    [2]RIKEN

Figure 1. Examples of label-noise robust conditional image generation. Each column shows samples belonging to the same class. In (c) and (d), each row contains samples generated with a fixed $z$ and a varied $y^g$. Our goal is, given *noisy* labeled data (b), to learn a conditional generative distribution that corresponds with *clean* labeled data (a). When naive cGAN (c) is trained with (b), it fails to learn the disentangled representations, disturbed by *noisy* labeled data. In contrast, proposed rGAN (d) succeeds in learning the representations disentangled on the basis of *clean* labels, which are close to (a), even when we can only access the *noisy* labeled data (b) during training.

## Abstract

*Generative adversarial networks (GANs) are a framework that learns a generative distribution through adversarial training. Recently, their class-conditional extensions (e.g., conditional GAN (cGAN) and auxiliary classifier GAN (AC-GAN)) have attracted much attention owing to their ability to learn the disentangled representations and to improve the training stability. However, their training requires the availability of large-scale accurate class-labeled data, which are often laborious or impractical to collect in a real-world scenario. To remedy this, we propose a novel family of GANs called label-noise robust GANs (rGANs), which, by incorporating a noise transition model, can learn a clean label conditional generative distribution even when training labels are noisy. In particular, we propose two variants: rAC-GAN, which is a bridging model between AC-GAN and the label-noise robust classification model, and rcGAN, which is an extension of cGAN and solves this problem with no reliance on any classifier. In addition to providing the theoretical background, we demonstrate the effectiveness of our models through extensive experiments using diverse GAN configurations, various noise settings, and multiple evaluation metrics (in which we tested 402 conditions in total). Our code is available at* `https://github.com/takuhirok/rGAN/`.

## 1. Introduction

In computer vision and machine learning, generative modeling has been actively studied to generate or reproduce samples indistinguishable from real data. Recently, deep generative models have emerged as a powerful framework for addressing this problem. Among them, generative adversarial networks (GANs) [16], which learn a generative distribution through adversarial training, have become a prominent one owing to their ability to learn any data distribution without explicit density estimation. This mitigates oversmoothing resulting from data distribution approximation, and GANs have succeeded in producing high-fidelity data for various tasks [30, 48, 82, 7, 25, 36, 64, 31, 79, 90, 23, 39, 10, 73, 72, 89, 8, 28].

Along with this success, various extensions of GANs have been proposed. Among them, class-conditional extensions (e.g., conditional GAN (cGAN) [47, 49] and auxiliary classifier GAN (AC-GAN) [52]) have attracted much attention mainly for two reasons. (1) By incorporating class labels as supervision, they can learn the representations that are disentangled between the class labels and the other factors. This allows them to selectively generate images conditioned on the class labels [47, 52, 28, 86, 29, 10]. Recently, this usefulness has also been demonstrated in class-specific data augmentation [14, 88]. (2) The added supervision sim-

plifies the learned target from an overall distribution to the conditional distribution. This helps stabilize the GAN training, which is typically unstable, and improves image quality [52, 49, 82, 7].

In contrast to these powerful properties, a possible limitation is that typical models rely on the availability of large-scale accurate class-labeled data and their performance depends on their accuracy. Indeed, as shown in Figure 1(c), when conventional cGAN is applied to noisy labeled data (where half labels are randomly flipped, as shown in Figure 1(b)), its performance is significantly degraded, influenced by the noisy labels. When datasets are constructed in real-world scenarios (e.g., crawled from websites or annotated via crowdsourcing), they tend to contain many mislabeled data (e.g., in Clothing1M [75], the overall annotation accuracy is only 61.54%). Therefore, this limitation would restrict application.

Motivated by these backgrounds, we address the following problem: *"How can we learn a clean label conditional distribution even when training labels are noisy?"* To solve this problem, we propose a novel family of GANs called *label-noise robust GANs* (*rGANs*) that incorporate a noise transition model representing a transition probability between the clean and noisy labels. In particular, we propose two variants: *rAC-GAN*, which is a bridging model between AC-GAN [52] and the label-noise robust classification model, and *rcGAN*, which is an extension of cGAN [47, 49] and solves this problem with no reliance on any classifier. As examples, we show generated image samples using rcGAN in Figure 1(d). As shown in this figure, our rcGAN is able to generate images conditioned on clean labels even where conventional cGAN suffers from severe degradation.

Another important issue regarding learning deep neural networks (DNNs) using noisy labeled data is the memorization effect. In image classification, a recent study [80] empirically demonstrated that DNNs can fit even noisy (or random) labels. Another study [5] experimentally showed that there are qualitative differences between DNNs trained on clean and noisy labeled data. To the best of our knowledge, no previous studies have sufficiently examined such an effect for conditional deep generative models. Motivated by these facts, in addition to providing a theoretical background on rAC-GAN and rcGAN, we conducted extensive experiments to examine the gap between theory and practice. In particular, we evaluated our models using diverse GAN configurations from standard to state-of-the-art in various label-noise settings including synthetic and real-world noise. We also tested our methods in the case when a noise transition model is known and in the case when it is not. Furthermore, we introduce an improved technique to stabilize training in a severely noisy setting (e.g., that in which 90% of the labels are corrupted) and show the effectiveness.

Overall, our contributions are summarized as follows:

- We tackle a novel problem called *label-noise robust conditional image generation*, in which the goal is to learn a clean label conditional generative distribution even when training labels are noisy.
- To solve this problem, we propose a new family of GANs called *rGANs* that incorporate a noise transition model into conditional extensions of GANs. In particular, we propose two variants, i.e., *rAC-GAN* and *rcGAN*, for the two representative class-conditional GANs, i.e., AC-GAN and cGAN.
- In addition to providing a theoretical background, we examine the gap between theory and practice through extensive experiments (in which we tested 402 conditions in total). Our code is available at https://github.com/takuhirok/rGAN/.

## 2. Related work

**Deep generative models.** Generative modeling has been a fundamental problem and has been actively studied in computer vision and machine learning. Recently, deep generative models have emerged as a powerful framework. Among them, three popular approaches are GANs [16], variational autoencoders (VAEs) [34, 61], and autoregressive models (ARs) [69]. All these models have pros and cons. One well-known problem with GANs is training instability; however, the recent studies have been making a great stride in solving this problem [12, 54, 62, 87, 3, 4, 44, 17, 30, 74, 48, 46, 82, 7]. In this paper, we focus on GANs because they have flexibility to the data representation, allowing for incorporating a noise transition model. However, with regard to VAEs and ARs, conditional extensions [33, 43, 78, 70, 58] have been proposed, and incorporating our ideas into them is a possible direction of future work.

**Conditional extensions of GANs.** As discussed in Section 1, conditional extensions of GANs have been actively studied to learn the representations that are disentangled between the conditional information and the other factors or to stabilize training and boost image quality. Other than class or attribute labels [47, 52, 28, 86, 29, 10], texts [56, 84, 83, 77], object locations [55], images [12, 25, 36, 73], or videos [72] are used as conditional information, and the effectiveness of conditional extensions of GANs has also been verified for them. In this paper, we focus on the situation in which noise exists in the label domain because obtaining robustness in such a domain has been a fundamental and important problem in image classification and has been actively studied, as discussed in the next paragraph. However, also in other domains (e.g., texts or images), it is highly likely that noise may exist when data are collected in real-world scenarios (e.g., crawled from websites or annotated via crowdsourcing). We believe that our findings would help the research also in these domains.

**Label-noise robust models.** Learning with noisy labels has been keenly studied since addressed in the learning theory community [1, 51]. Lately, this problem has also been studied in image classification with DNNs. For instance, to obtain label-noise robustness, one approach replaces a typical cross-entropy loss with a noise-tolerant loss [2, 85]. Another approach cleans up labels or selects clean labels out of noisy labels using neural network predictions or gradient directions [57, 66, 42, 26, 60, 19]. The other approach incorporates a noise transition model [65, 27, 53, 15], similarly to ours. These studies show promising results in both theory and practice and our study is based on their findings.

The main difference from them is that their goal is to obtain label-noise robustness in image classification, but our goal is to obtain such robustness in conditional image generation. We remark that our developed rAC-GAN internally uses a classifier; thus, it can be viewed as a bridging model between noise robust image classification and conditional image generation. Note that we also developed rcGAN, which is a classifier-free model, motivated by the recent studies [52, 49] that indicate that AC-GAN tends to lose diversity through a side effect of generating recognizable (i.e., classifiable) images. Another related topic is *pixel*-noise robust image generation [6, 37]. The difference from them is that they focused on the noise inserted in a *pixel* domain, but we focus on the noise in a *label* domain.

## 3. Notation and problem statement

We begin by defining notation and the problem statement. Throughout, we use superscript $r$ to denote the real distribution and $g$ the generative distribution. Let $\boldsymbol{x} \in \mathcal{X}$ be the target data (e.g., images) and $y \in \mathcal{Y}$ the corresponding class label. Here, $\mathcal{X}$ is the data space $\mathcal{X} \subseteq \mathbb{R}^d$, where $d$ is the dimension of the data, and $\mathcal{Y}$ is the label space $\mathcal{Y} = \{1, \ldots, c\}$, where $c$ is the number of classes. We assume that $y$ is noisy (and we denote such noisy label by $\tilde{y}$) and there exists a corresponding clean label $\hat{y}$ that we cannot observe during training. In particular, we assume *class-dependent* noise in which each clean label $\hat{y} = i$ is corrupted to a noisy label $\tilde{y} = j$ with a probability $p(\tilde{y} = j | \hat{y} = i) = T_{i,j}$, independently of $\boldsymbol{x}$, where we define a noise transition matrix as $T = (T_{i,j}) \in [0, 1]^{c \times c}$ ($\sum_i T_{i,j} = 1$). Note that this assumption is commonly used in label-noise robust image classification (e.g., [2, 85, 65, 27, 53, 15]).

Our task is, when given noisy labeled samples $(\boldsymbol{x}^r, \tilde{y}^r) \sim \tilde{p}^r(\boldsymbol{x}, \tilde{y})$, to construct a label-noise robust conditional generator such that $\hat{p}^g(\boldsymbol{x}, \hat{y}) = \hat{p}^r(\boldsymbol{x}, \hat{y})$, which can generate $\boldsymbol{x}$ conditioned on *clean* $\hat{y}$ rather than conditioned on *noisy* $\tilde{y}$. This task is challenging for typical conditional generative models, such as AC-GAN [52] (Figure 2(b)) and cGAN [47, 49] (Figure 2(d)), because they attempt to construct a generator conditioned on the observable labels; i.e., in this case, they attempt to construct a *noisy*-label-

dependent generator that generates $\boldsymbol{x}$ conditioned on *noisy* $\tilde{y}$ rather than conditioned on *clean* $\hat{y}$. Our main idea for solving this problem is to incorporate a noise transition model, i.e., $p(\tilde{y}|\hat{y})$, into these models (viewed as orange rectangles in Figures 2(c) and (e)). In particular, we develop two variants: rAC-GAN and rcGAN. We describe their details in Sections 4 and 5, respectively.

## 4. Label-noise robust AC-GAN: rAC-GAN

### 4.1. Background: AC-GAN

AC-GAN [52] is one of representative conditional extensions of GANs [16]. AC-GAN learns a conditional generator $G$ that transforms noise $\boldsymbol{z}$ and label $y^g$ into data $\boldsymbol{x}^g = G(\boldsymbol{z}, y^g)$ with two networks. One is a discriminator $D$ that assigns probability $p = D(\boldsymbol{x})$ for samples $\boldsymbol{x} \sim p^r(\boldsymbol{x})$ and assigns $1 - p$ for samples $\boldsymbol{x} \sim p^g(\boldsymbol{x})$. The other is an auxiliary classifier $C(y|\boldsymbol{x})$ that represents a probability distribution over class labels given $\boldsymbol{x}$. These networks are optimized by using two losses, namely, an adversarial loss and an auxiliary classifier loss.

**Adversarial loss.** An adversarial loss is defined as

$$
\begin{aligned}
\mathcal{L}_{\mathrm{GAN}} = \; & \mathbb{E}_{\boldsymbol{x}^r \sim p^r(\boldsymbol{x})}[\log D(\boldsymbol{x}^r)] \\
& + \mathbb{E}_{\boldsymbol{z} \sim p(\boldsymbol{z}), y^g \sim p(y)}[\log(1 - D(G(\boldsymbol{z}, y^g)))], \quad (1)
\end{aligned}
$$

where $D$ attempts to find the best decision boundary between real and generated data by maximizing this loss, and $G$ attempts to generate data indistinguishable by $D$ by minimizing this loss.

**Auxiliary classifier loss.** An auxiliary classifier loss is used to make the generated data belong to the target class. To achieve this, first $C$ is optimized using a classification loss of real data:

$$
\mathcal{L}_{\mathrm{AC}}^r = \mathbb{E}_{(\boldsymbol{x}^r, y^r) \sim p^r(\boldsymbol{x}, y)}[-\log C(y = y^r | \boldsymbol{x}^r)], \quad (2)
$$

where $C$ learns to classify real data to the corresponding class by minimizing this loss. Then, $G$ is optimized by using a classification loss of generated data:

$$
\mathcal{L}_{\mathrm{AC}}^g = \mathbb{E}_{\boldsymbol{z} \sim p(\boldsymbol{z}), y^g \sim p(y)}[-\log C(y = y^g | G(\boldsymbol{z}, y^g))], \quad (3)
$$

where $G$ attempts to generate data belonging to the corresponding class by minimizing this loss.

**Full objective.** In practice, shared networks between $D$ and $C$ are commonly used [52, 17]. In this setting, the full objective is written as

$$
\begin{aligned}
\mathcal{L}_{D/C} &= -\mathcal{L}_{\mathrm{GAN}} + \lambda_{\mathrm{AC}}^r \mathcal{L}_{\mathrm{AC}}^r, && (4) \\
\mathcal{L}_G &= \mathcal{L}_{\mathrm{GAN}} + \lambda_{\mathrm{AC}}^g \mathcal{L}_{\mathrm{AC}}^g, && (5)
\end{aligned}
$$

where $\lambda_{\mathrm{AC}}^r$ and $\lambda_{\mathrm{AC}}^g$ are the trade-off parameters between the adversarial loss and the auxiliary classifier loss for the real and generated data, respectively. $D/C$ and $G$ are optimized by minimizing $\mathcal{L}_{D/C}$ and $\mathcal{L}_G$, respectively.
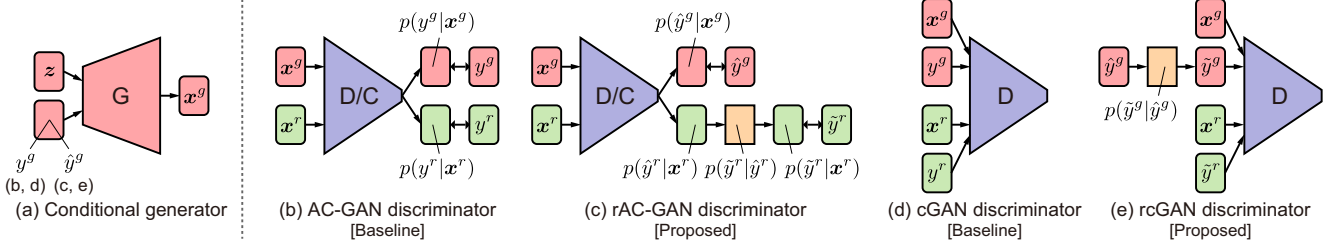
Figure 2. Comparison of naive and label-noise robust GANs. We denote the generator, discriminator, and auxiliary classifier by $G$, $D$, and $C$, respectively. Among all models, conditional generators (a) are similar. In our rAC-GAN (c) and rcGAN (e), we incorporate a noise transition model (viewed as an orange rectangle) into AC-GAN (b) and cGAN (d), respectively.

## 4.2. rAC-GAN

By the above definition, when $y^r$ is noisy (i.e., $\tilde{y}^r$ is given) and $C$ fits such noisy labels,[1] AC-GAN learns the *noisy* label conditional generator $G(\boldsymbol{z}, \tilde{y}^g)$. In contrast, our goal is to construct the *clean* label conditional generator $G(\boldsymbol{z}, \hat{y}^g)$. To achieve this goal, we incorporate a noise transition model (i.e., $p(\tilde{y}|\hat{y})$; viewed as an orange rectangle in Figure 2(c)) into the auxiliary classifier. In particular, we reformulate the auxiliary classifier loss as

$$\mathcal{L}^r_{\mathrm{rAC}} = \mathbb{E}_{(\boldsymbol{x}^r, \tilde{y}^r) \sim \tilde{p}^r(\boldsymbol{x}, \tilde{y})}[-\log \tilde{C}(\tilde{y} = \tilde{y}^r | \boldsymbol{x}^r)]$$
$$= \mathbb{E}_{(\boldsymbol{x}^r, \tilde{y}^r) \sim \tilde{p}^r(\boldsymbol{x}, \tilde{y})}$$
$$[-\log \sum_{\hat{y}^r} p(\tilde{y} = \tilde{y}^r | \hat{y} = \hat{y}^r)\hat{C}(\hat{y} = \hat{y}^r | \boldsymbol{x}^r)]$$
$$= \mathbb{E}_{(\boldsymbol{x}^r, \tilde{y}^r) \sim \tilde{p}^r(\boldsymbol{x}, \tilde{y})}[-\log \sum_{\hat{y}^r} T_{\hat{y}^r, \tilde{y}^r}\hat{C}(\hat{y} = \hat{y}^r | \boldsymbol{x}^r)], \quad (6)$$

where we denote the *noisy* label classifier by $\tilde{C}$ and the *clean* label classifier by $\hat{C}$ (and we explain the reason why we call it *clean* in Theorem 1). Between the first and second lines, we assume that the noise transition is independent of $\boldsymbol{x}$, as mentioned in Section 3. Note that this formulation (called the *forward correction*) is often used in label-noise robust classification models [65, 27, 53, 15] and rAC-GAN can be viewed as a bridging model between GANs and them. In naive AC-GAN, $\tilde{C}$ is optimized for $\mathcal{L}^r_{\mathrm{AC}}$, whereas in our rAC-GAN, $\hat{C}$ is optimized for $\mathcal{L}^r_{\mathrm{rAC}}$. Similarly, $G$ is optimized using $\hat{C}$ rather than using $\tilde{C}$:

$$\mathcal{L}^g_{\mathrm{rAC}} = \mathbb{E}_{\boldsymbol{z} \sim p(\boldsymbol{z}), \hat{y}^g \sim p(\hat{y})}[-\log \hat{C}(\hat{y} = \hat{y}^g | G(\boldsymbol{z}, \hat{y}^g))]. \quad (7)$$

**Theoretical background.** In the above, we use a cross-entropy loss, which is a kind of proper composite loss [59]. In this case, Theorem 2 in [53] shows that minimizing the

---

[1]Zhang et al. [80] discuss generalization and memorization of DNNs and empirically demonstrated that DNNs are capable of fitting even noisy (or random) labels. Although other studies empirically demonstrated that some techniques (e.g., dropout [5], mixup [81], and high learning rate [66]) are useful for preventing DNNs from memorizing noisy labels, their theoretical support still remains as an open issue. In this paper, we conducted experiments on various GAN configurations to investigate such effect in our task. See Section 7.1 for details.

forward corrected loss (i.e., Equation 6) is equal to minimizing the original loss under the clean distribution. More precisely, the following theorem holds.

**Theorem 1.** *When $T$ is nonsingular,*

$$\operatorname*{argmin}_{\hat{C}} \mathbb{E}_{(\boldsymbol{x}^r, \tilde{y}^r) \sim \tilde{p}^r(\boldsymbol{x}, \tilde{y})}[-\log \sum_{\hat{y}^r} T_{\hat{y}^r, \tilde{y}^r}\hat{C}(\hat{y} = \hat{y}^r | \boldsymbol{x}^r)]$$
$$= \operatorname*{argmin}_{\hat{C}} \mathbb{E}_{(\boldsymbol{x}^r, \hat{y}^r) \sim \hat{p}^r(\boldsymbol{x}, \hat{y})}[-\log \hat{C}(\hat{y} = \hat{y}^r | \boldsymbol{x}^r)]. \quad (8)$$

For a detailed proof, refer to Theorem 2 in [53]. This supports the idea that, by minimizing $\mathcal{L}^r_{\mathrm{rAC}}$ for noisy labeled samples, we can obtain $\hat{C}$ that classifies $\boldsymbol{x}$ as its corresponding clean label $\hat{y}$. In rAC-GAN, $G$ is optimized for this *clean* classifier $\hat{C}$; hence, in $G$'s input space, $\hat{y}^g$ is encouraged to represent clean labels.

## 5. Label-noise robust cGAN: rcGAN

### 5.1. Background: cGAN

cGAN [47, 49] is another representative conditional extension of GANs [16]. In cGAN, a conditional generator $G(\boldsymbol{z}, y^g)$ and a conditional discriminator $D(\boldsymbol{x}, y)$ are jointly trained using a conditional adversarial loss.

**Conditional adversarial loss.** A conditional adversarial loss is defined as

$$\mathcal{L}_{\mathrm{cGAN}} = \mathbb{E}_{(\boldsymbol{x}^r, y^r) \sim p^r(\boldsymbol{x}, y)}[\log D(\boldsymbol{x}^r, y^r)]$$
$$+ \mathbb{E}_{\boldsymbol{z} \sim p(\boldsymbol{z}), y^g \sim p(y)}[\log(1 - D(G(\boldsymbol{z}, y^g), y^g))], \quad (9)$$

where $D$ attempts to find the best decision boundary between real and generated data conditioned on $y$ by maximizing this loss. In contrast, $G$ attempts to generate data indistinguishable by $D$ with a constraint on $y^g$ by minimizing this loss. In an optimal condition [16], cGAN learns $G(\boldsymbol{z}, y)$ such that $p^g(\boldsymbol{x}, y) = p^r(\boldsymbol{x}, y)$.

### 5.2. rcGAN

By the above definition, when $y^r$ is noisy (i.e., $\tilde{y}^r$ is given), cGAN learns the *noisy* label conditional generator $G(\boldsymbol{z}, \tilde{y}^g)$. In contrast, our goal is to construct the *clean* label conditional generator $G(\boldsymbol{z}, \hat{y}^g)$. To achieve this goal, we

insert a noise transition model (viewed as an orange rectangle in Figure 2(e)) before $\hat{y}^g$ is given to $D$. In particular, we sample $\tilde{y}^g$ from $\tilde{y}^g \sim p(\tilde{y}|\hat{y}^g)$ and redefine Equation 9 as

$$\mathcal{L}_{\text{rcGAN}} = \mathbb{E}_{(\boldsymbol{x}^r, \tilde{y}^r) \sim \tilde{p}^r(\boldsymbol{x}, \tilde{y})}[\log D(\boldsymbol{x}^r, \tilde{y}^r)] \\ + \mathbb{E}_{\boldsymbol{z} \sim p(\boldsymbol{z}), \hat{y}^g \sim p(\hat{y}), \tilde{y}^g \sim p(\tilde{y}|\hat{y}^g)}[\log(1 - D(G(\boldsymbol{z}, \hat{y}^g), \tilde{y}^g))], \tag{10}$$

where $D$ attempts to find the best decision boundary between real and generated data conditioned on *noisy* labels $\tilde{y}$, by maximizing this loss. In contrast, $G$ attempts to generate data indistinguishable by $D$ with a constraint on *clean* labels $\hat{y}^g$ (and we explain the rationale behind calling it *clean* in Theorem 2), by minimizing this loss.

**Theoretical background.** In an optimal condition, the following theorem holds.

**Theorem 2.** *When $T$ is nonsingular (i.e., $T$ has a unique inverse), $G$ is optimal if and only if $\hat{p}^g(\boldsymbol{x}, \hat{y}) = \hat{p}^r(\boldsymbol{x}, \hat{y})$.*

*Proof.* For $G$ fixed, rcGAN is the same as cGAN where $y$ is replaced by $\tilde{y}$. Therefore, by extending Proposition 1 and Theorem 1 in [16] (GAN optimal solution) to a conditional setting, the optimal discriminator $D$ for fixed $G$ is

$$D(\boldsymbol{x}, \tilde{y}) = \frac{\tilde{p}^r(\boldsymbol{x}, \tilde{y})}{\tilde{p}^r(\boldsymbol{x}, \tilde{y}) + \tilde{p}^g(\boldsymbol{x}, \tilde{y})}. \tag{11}$$

Then $G$ is optimal if and only if

$$\tilde{p}^g(\boldsymbol{x}, \tilde{y}) = \tilde{p}^r(\boldsymbol{x}, \tilde{y}). \tag{12}$$

As mentioned in Section 3, we assume that label corruption occurs with $p(\tilde{y}|\hat{y})$, i.e., independently of $\boldsymbol{x}$. In this case,

$$\tilde{p}(\boldsymbol{x}, \tilde{y}) = \tilde{p}(\tilde{y}|\boldsymbol{x})p(\boldsymbol{x}) = \sum_{\hat{y}} p(\tilde{y}|\hat{y})\hat{p}(\hat{y}|\boldsymbol{x})p(\boldsymbol{x}) \\ = \sum_{\hat{y}} p(\tilde{y}|\hat{y})\hat{p}(\boldsymbol{x}, \hat{y}) = \sum_{\hat{y}} T_{\hat{y}, \tilde{y}}\hat{p}(\boldsymbol{x}, \hat{y}). \tag{13}$$

Substituting Equation 13 into Equation 12 gives

$$\sum_{\hat{y}} T_{\hat{y}, \tilde{y}}\hat{p}^g(\boldsymbol{x}, \hat{y}) = \sum_{\hat{y}} T_{\hat{y}, \tilde{y}}\hat{p}^r(\boldsymbol{x}, \hat{y}). \tag{14}$$

By considering the matrix form,

$$T^\top \hat{P}^g = T^\top \hat{P}^r, \tag{15}$$

where $\hat{P}^g = [\hat{p}^g(\boldsymbol{x}, \hat{y} = 1), \ldots, \hat{p}^g(\boldsymbol{x}, \hat{y} = c)]^\top$ and $\hat{P}^r = [\hat{p}^r(\boldsymbol{x}, \hat{y} = 1), \ldots, \hat{p}^r(\boldsymbol{x}, \hat{y} = c)]^\top$. When $T$ has an inverse,

$$T^\top \hat{P}^g = T^\top \hat{P}^r \Leftrightarrow \hat{P}^g = (T^\top)^{-1}T^\top \hat{P}^r = \hat{P}^r. \tag{16}$$

As the corresponding elements in $\hat{P}^g$ and $\hat{P}^r$ are equal, $\hat{p}^g(\boldsymbol{x}, \hat{y}) = \hat{p}^r(\boldsymbol{x}, \hat{y})$. □

This supports the idea that, in an optimal condition, rcGAN learns $G(\boldsymbol{z}, \hat{y})$ such that $\hat{p}^g(\boldsymbol{x}, \hat{y}) = \hat{p}^r(\boldsymbol{x}, \hat{y})$.

## 6. Advanced techniques for practice

### 6.1. Noise transition probability estimation

In the above, we assume that $T$ is known, but this assumption may be too strict for real-world applications. However, fortunately, previous studies [65, 27, 53, 15] have been eagerly tackling this problem and several methods for estimating $T'$ (where we denote the estimated $T$ by $T'$) have been proposed. Among them, we tested a *robust two-stage training algorithm* [53] in the experiments and analyzed the effects of estimated $T'$. We show the results in Section 7.2.

### 6.2. Improved technique for severely noisy data

Thorough extensive experiments, we find that some GAN configurations suffer from performance degradation in a severely noisy setting (e.g., in which 90% of the labels are corrupted). In this type of environment, each label is flipped with a high probability. This disturbs $G$ form associating an image with a label. To strengthen their connection, we incorporate mutual information regularization [9]:

$$\mathcal{L}_{\text{MI}} = \mathbb{E}_{\boldsymbol{z} \sim p(\boldsymbol{z}), \hat{y}^g \sim p(\hat{y})}[-\log Q(\hat{y} = \hat{y}^g|G(\boldsymbol{z}, \hat{y}^g))], \tag{17}$$

where $Q(\hat{y}|\boldsymbol{x})$ is an auxiliary distribution approximating a true posterior $p(\hat{y}|\boldsymbol{x})$. We optimize $G$ and $Q$ by minimizing this loss with trade-off parameters $\lambda_{\text{MI}}^g$ and $\lambda_{\text{MI}}^q$, respectively. This formulation is similar to Equation 7, but the difference is whether $G$ is optimized for $\hat{C}$ (optimized using real images and noisy labels) or for $Q$ (optimized using generated images and clean labels). We demonstrate the effectiveness of this technique in Section 7.3.

## 7. Experiments

### 7.1. Comprehensive study

In Sections 4 and 5, we showed that our approach is theoretically grounded. However, generally, in DNNs, there is still a gap between theory and practice. In particular, the label-noise effect in DNNs just recently began to be discussed in image classification [80, 5], and it is demonstrated that such a gap exists. However, in conditional image generation, such an effect has not been sufficiently examined. To advance this research, we first conducted a comprehensive study, i.e., compared the performance of conventional AC-GAN and cGAN and proposed rAC-GAN and rcGAN using diverse GAN configurations in various label-noise settings with multiple evaluation metrics.[2] Due to the space limitation, we briefly review the experimental setup and only provide the important results in this main text. See the Appendix and our website for details and more results.

**Dataset.** We verified the effectiveness of our method on two benchmark datasets: CIFAR-10 and CIFAR-100 [35],

---

[2]Through Sections 7.1–7.3, we tested 392 conditions in total. For each condition, we trained two models with different initializations and report the results averaged over them.

Figure 3. Generated image samples on CIFAR-10. Each column shows samples belonging to the same class. Each row contains samples generated with a fixed $z$ and a varied $y^g$. In symmetric noise (a), cSN-GAN is primarily influenced by noisy labels and fails to learn the disentangled representations. In asymmetric noise (b), it is expected that fourth and sixth columns will include cat and dog, respectively. However, in AC-CT-GAN and cSN-GAN, these columns contain the inverse. As evidence, we list the accuracy in the fourth column for cat/dog classes in Table 1. These scores indicate that the proposed models are robust but the baselines are weak for the flipped classes. See Figures 6–9 in the Appendix for more samples.

which are commonly used in both image generation and label-noise robust image classification. Both datasets contain $60k$ $32 \times 32$ natural images, which are divided into $50k$ training and $10k$ test images. CIFAR-10 and CIFAR-100 have 10 and 100 classes, respectively. We assumed two label-noise settings that are popularly used in label-noise robust image classification: (1) **Symmetric** (class-independent) noise [71]: For all classes, ground truth labels are replaced with uniform random classes with probability $\mu$. (2) **Asymmetric** (class-dependent) noise [53]: Ground truth labels are flipped with probability $\mu$ by mimicking real mistakes between similar classes. Following [53], for CIFAR-10, ground truth labels are replaced with *truck* $\rightarrow$ *automobile*, *bird* $\rightarrow$ *airplane*, *deer* $\rightarrow$ *horse*, and *cat* $\leftrightarrow$ *dog*, and for CIFAR-100, ground truth labels are flipped into the next class circularly within the same superclasses. In both settings, we selected $\mu$ from $\{0.1, 0.3, 0.5, 0.7, 0.9\}$.

**GAN configurations.** A recent study [40] shows the sensitivity of GANs to hyperparameters. However, when clean labeled data are not available, it is impractical to tune the hyperparameters for each label-noise setting. Hence, instead of searching for the best model with hyperparameter tuning, we tested various GAN configurations using the default parameters that are typically used in clean label settings and examined the label-noise effect. We chose four models to cover standard, widely accepted, and state-of-the-art models: **DCGAN** [54], **WGAN-GP** [17], **CT-GAN** [74], and **SN-GAN** [48]. We implemented AC-GAN, rAC-GAN, cGAN, and rcGAN based on them. For cGAN and rcGAN, we used the *concat* discriminator [47] for DCGAN and the *projection* discriminator [49] for the others.

**Evaluation metrics.** As discussed in previous studies [67, 40, 63], evaluation and comparison of GANs can be challenging partially because of the lack of an explicit likelihood measure. Considering this fact, we used four metrics for a comprehensive analysis: (1) the Fréchet Inception distance (**FID**), (2) **Intra FID**, (3) the **GAN-test**, and (4) the **GAN-train**. The FID [22] measures the distance between

| | AC-CT-GAN | rAC-CT-GAN | cSN-GAN | rcSN-GAN |
|---|---|---|---|---|
| cat/dog | 13.4/**83.9** | **84.8**/10.3 | 35.6/**55.9** | **75.9**/13.0 |

Table 1. Accuracy in the fourth column in Figure 3(b) (ground truth: cat) for the flipped classes (cat $\leftrightarrow$ dog)

$p^r$ and $p^g$ in Inception embeddings. We used it to assess the quality of an overall generative distribution. Intra FID [49] calculates the FID for each class. We used it to assess the quality of a conditional generative distribution.[3] The GAN-test [63] is the accuracy of a classifier trained on real images and evaluated on generated images. This metric approximates the precision (image quality) of GANs. The GAN-train [63] is the accuracy of a classifier trained on generated images and evaluated on real images in a test. This metric approximates the recall (diversity) of GANs.

**Results.** We present the quantitative results for each condition in Figure 4 and provide a comparative summary between the proposed models (i.e., rAC-GAN and rcGAN) and the baselines (i.e., AC-GAN and cGAN) across all conditions in Figure 5. We show the samples of generated images on CIFAR-10 with $\mu = 0.7$ in Figure 3. Regarding the FID (i.e., evaluating the quality of the overall generative distribution), the baselines and the proposed models are comparable in most cases, but when we use CT-GAN and SN-GAN (i.e., state-of-the-art models) in symmetric noise, the proposed models tend to outperform the baselines (32/40 conditions). This indicates that the label ambiguity caused by symmetric noise could disturb the learning of GANs if they have the high data-fitting ability. However, this degradation can be mitigated by using the proposed methods.

Regarding the other metrics (i.e., evaluating the quality of the conditional generative distribution), rAC-GAN and rcGAN tend to outperform AC-GAN and cGAN, respectively, across all the conditions. The one exception is rAC-WGAN-GP on CIFAR-10 with symmetric noise, but we find that it can be improved using the technique introduced in Section 6.2. We demonstrate this in Section 7.3. Among the four models, CT-GAN and SN-GAN work relatively

---

[3]We used Intra FID only for CIFAR-10 because, in CIFAR-100, the number of clean labeled data for each class (500) is insufficient.
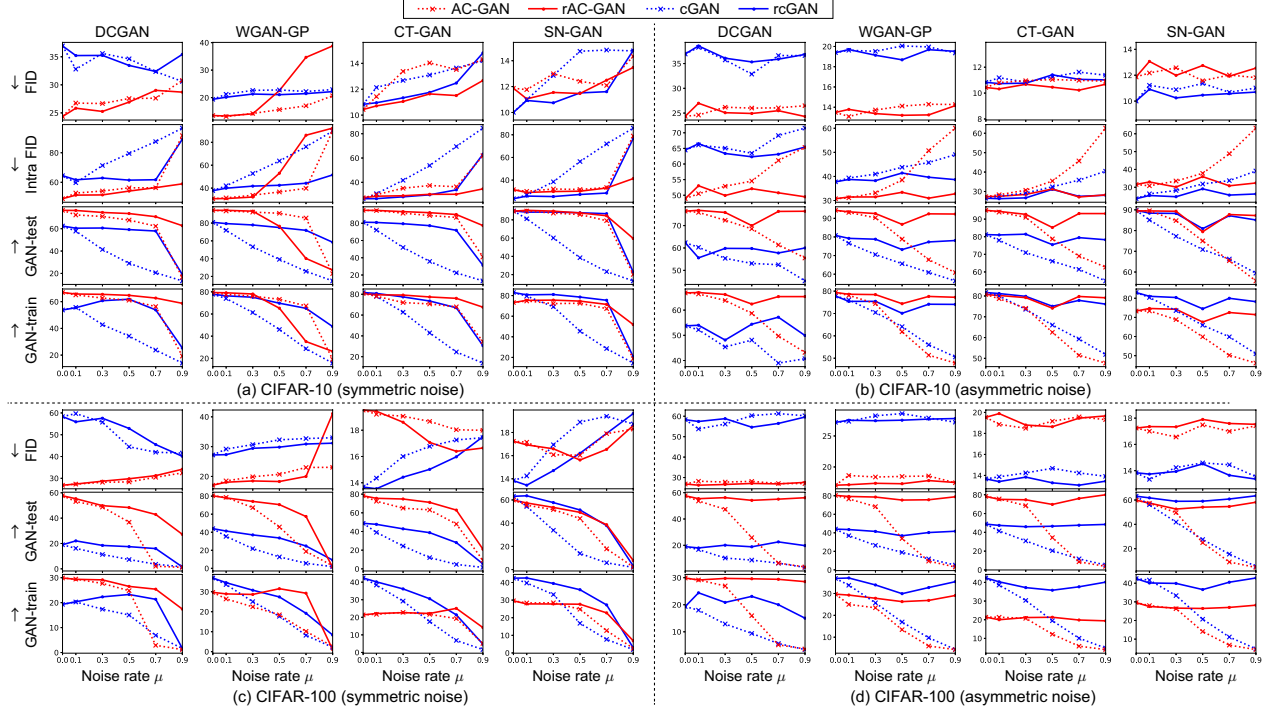
Figure 4. Quantitative results on CIFAR-10 and CIFAR-100. ↓ indicates the smaller the value, the better the performance. ↑ indicates the larger the value, the better the performance. Note that the scale is adjusted on each graph for easy viewing.
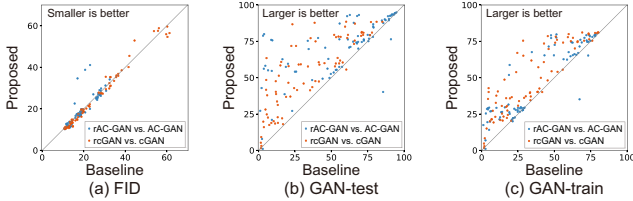


Figure 5. Comparison between the proposed models and the baselines across all the conditions in Figure 4.

|  | AC-GAN | rAC-GAN | cGAN | rcGAN |
|---|---|---|---|---|
| Symmetric | -0.846 ± 0.084 | -0.786 ± 0.163 | **-0.989** ± 0.013 | -0.818 ± 0.142 |
| Asymmetric | **-0.976** ± 0.008 | -0.476 ± 0.119 | **-0.985** ± 0.029 | -0.427 ± 0.274 |

Table 2. Pearson correlation coefficient between the noise rate and GAN-train. The scores are averaged over all GAN configurations.

well for rAC-GAN and rcGAN, respectively. This tendency is also observed in clean label settings (i.e., $\mu = 0$). This indicates that the performance of rAC-GAN and rcGAN is closely related to the advance in the baseline GANs.

An interesting finding is that the performance of cGAN in Intra FID, GAN-test, and GAN-train degrades linearly depending on the noise rate. To confirm this numerically, we calculated the Pearson correlation coefficient between the GAN-train and the noise rate. We list these in Table 2. These scores confirm that cGAN has the highest dependency on the noise rate, i.e., cGAN can fit even nosy labels. In contrast, AC-GAN shows robustness for symmetric noise but weakness for asymmetric noise. This would be related to the difficulty of memorization. In symmetric noise, the corruption variety is large, making it difficult to memorize labels. As a result, AC-GAN prioritizes learning simple (i.e., clean) labels, in a similar way as DNNs in image clas-

sification [5]. In contrast, in asymmetric noise, the label corruption pattern is restrictive; as a result, AC-GAN easily fits noisy labels. Unlike AC-GAN, cGAN is a classifier-free model; therefore, cGAN tends to fit the given labels regardless of whether labels are noisy or not.

## 7.2. Effects of estimated $T'$

In Section 7.1, we report the results using known $T$. As a more practical setting, we also evaluate our method with $T'$ estimated by a robust two-stage training algorithm [53]. We used CT-GAN for rAC-GAN and SN-GAN for rcGAN, which worked relatively well in both noisy and clean settings in Section 7.1. We list the scores in Table 3. In CIFAR-10, even using $T'$, rAC-CT-GAN and rcSN-GAN tend to outperform conventional AC-CT-GAN and cSN-GAN, respectively, and show robustness to label noise. In CIFAR-100, when the noise rate is low, rAC-CT-GAN and rcSN-GAN work moderately well; however, in highly noisy settings, their performance is degraded. Note that such a tendency has also been observed in noisy label image classification with $T'$ [53], in which the authors argue that the high-rate mixture and limited number of images per class (500) make it difficult to estimate the correct $T$. Further improvement remains as an open issue.

## 7.3. Evaluation of improved technique

As shown in Figure 4, rAC-GAN and rcGAN show robustness for label noise in almost all cases, but we find that they are still weak to severely noisy settings (i.e., symmet-

| Model | Metric | CIFAR-10 (symmetric noise) | | | | | CIFAR-10 (asymmetric noise) | | | | | CIFAR-100 (symmetric noise) | | | | | CIFAR-100 (asymmetric noise) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 0.1 | 0.3 | 0.5 | 0.7 | 0.9 | 0.1 | 0.3 | 0.5 | 0.7 | 0.9 | 0.1 | 0.3 | 0.5 | 0.7 | 0.9 | 0.1 | 0.3 | 0.5 | 0.7 | 0.9 |
| rAC-CT-GAN with $T'$ | FID ↓ | 10.9 | 11.4 | 11.3 | 11.5 | 13.0 | 10.8 | 10.2 | 10.2 | 10.4 | 11.0 | 19.7 | 19.3 | 17.7 | 17.3 | 18.5 | 19.4 | 19.3 | 19.7 | 18.8 | 19.0 |
| | Intra FID ↓ | 28.7 | **31.0** | **30.1** | 31.7 | **38.9** | 28.5 | **27.4** | **31.2** | **35.0** | **36.8** | – | – | – | – | – | – | – | – | – | – |
| | GAN-test ↑ | 95.3 | 93.2 | **92.0** | 87.7 | **70.4** | 94.9 | 92.9 | **85.2** | 78.5 | 76.6 | **76.6** | 67.1 | **68.1** | *1.0* | 2.5 | 74.1 | 68.9 | **28.7** | 7.2 | 2.2 |
| | GAN-train ↑ | 78.7 | **75.9** | **76.9** | **73.7** | **63.4** | 79.8 | **79.5** | **74.0** | **69.1** | **67.3** | 21.2 | 21.4 | 23.3 | *1.0* | 2.3 | 19.1 | 19.9 | 10.7 | 5.5 | 3.9 |
| rcSN-GAN with $T'$ | FID ↓ | 10.7 | 11.9 | 12.4 | 12.1 | 15.0 | 10.8 | 10.8 | 11.0 | 10.9 | 11.3 | 14.3 | 16.6 | 17.5 | 20.0 | 19.8 | 13.8 | 14.1 | 14.7 | 14.7 | 13.9 |
| | Intra FID ↓ | 25.5 | **29.4** | **29.4** | 29.7 | 87.4 | 25.7 | 26.0 | **28.7** | 32.6 | 33.9 | – | – | – | – | – | – | – | – | – | – |
| | GAN-test ↑ | **85.3** | 79.0 | **84.8** | 82.8 | 15.9 | 86.6 | **87.2** | **84.0** | 74.9 | **71.2** | 53.4 | 36.6 | **37.7** | *1.0* | 1.7 | **65.0** | **63.0** | **32.4** | *7.8* | 3.8 |
| | GAN-train ↑ | 80.7 | **78.1** | **77.4** | **75.6** | 15.0 | 80.5 | **79.0** | **75.7** | **69.3** | **65.7** | 40.1 | 32.8 | **31.3** | *1.0* | 1.8 | **41.7** | **39.3** | 20.1 | *6.1* | 3.9 |

Table 3. Quantitative results using the estimated $T'$. The second row indicates a noise rate. Bold and italic fonts indicate that the score is better or worse by more than 3 points over or under the baseline models (i.e., AC-CT-GAN or cSN-GAN), respectively. See Table 6 and Figure 11 in the Appendix for more detailed comparison and generated image samples, respectively.

| Model | Metric | CIFAR-10 (symmetric noise) | | | | CIFAR-100 (symmetric noise) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | A | B | C | D | A | B | C | D |
| Improved rAC-GAN | FID ↓ | 27.9 | **14.7** | 12.4 | 13.5 | 33.1 | **20.4** | 17.2 | 18.4 |
| | Intra FID ↓ | **55.7** | **34.6** | 33.4 | **36.9** | – | – | – | – |
| | GAN-test ↑ | 65.1 | **77.7** | 78.2 | **63.5** | 26.2 | **22.5** | 21.5 | **15.4** |
| | GAN-train ↑ | 59.9 | **70.8** | 69.1 | **59.7** | 17.1 | **16.3** | 14.8 | **11.7** |
| Improved rcGAN | FID ↓ | 30.4 | **16.9** | 14.2 | 14.9 | *50.2* | **25.8** | 18.0 | 18.7 |
| | Intra FID ↓ | **76.9** | **39.6** | **52.9** | **48.2** | – | – | – | – |
| | GAN-test ↑ | 27.3 | **65.7** | 38.9 | **48.8** | 4.5 | **12.0** | 9.5 | **6.1** |
| | GAN-train ↑ | **31.9** | **60.7** | 36.7 | **47.3** | **6.0** | 10.3 | 7.5 | 4.4 |

Table 4. Quantitative results using the improved technique. In the second row, A, B, C, and D indicate DCGAN, WGAN-GP, CT-GAN, and SN-GAN, respectively. We evaluated in severely noisy settings (i.e., symmetric noise with $\mu = 0.9$). Bold and italic fonts indicate that the score is better or worse by more than 3 points over or under naive models (i.e., rAC-GAN or rcGAN), respectively. See Table 7 and Figure 12 in the Appendix for more detailed comparison and generated image samples, respectively.

| Metric | Clean | | Noisy | | | | Mixed | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | AC | c | AC | rAC | c | rc | AC | rAC | c | rc |
| FID ↓ | 6.8 | 12.0 | **4.4** | 4.6 | 9.4 | 9.4 | 4.8 | **4.7** | 10.5 | **9.7** |
| GAN-train ↑ | 56.6 | 53.9 | 49.5 | **51.7** | 48.6 | **49.8** | 52.8 | **57.0** | 51.7 | **55.0** |

Table 5. Quantitative results on Clothing1M. AC, rAC, c, and rc denote AC-CT-GAN, rAC-CT-GAN, cSN-GAN, and rcSN-GAN, respectively. Bold font indicates better scores in each block. See Figure 13 in the Appendix for generated image samples.

ric noise with $\mu = 0.9$) even though using known $T$. To improve the performance, we developed an improved technique (Section 6.2). In this section, we validate its effect. We list the scores in Table 4. We find that the improved degree depends on the GAN configurations, but, on the whole, the performance is improved by the proposed technique. In particular, we find that the improved technique is most effective for rAC-WGAN-GP, in which all the scores doubled compared to those of naive rAC-WGAN-GP.

### 7.4. Evaluation on real-world noise

Finally, we tested on Clothing1M [75] to analyze the effectiveness on real-world noise.[4] Clothing1M contains $1M$ clothing images in 14 classes. The data are collected from several online shopping websites and include many mislabeled samples. This dataset also contains $50k$, $14k$, and $10k$ of clean data for training, validation, and testing, respectively. Following the previous studies [75, 53], we approximated $T$ using the partial ($25k$) training data that have both clean and noisy labels. We tested on three settings: (1) $50k$ **clean** data, (2) $1M$ **noisy** data, and (3) **mixed** data

that consists of clean data (bootstrapped to $500k$) and $1M$ noisy data, which are used in [75] to boost the performance of image classification. We used AC-CT-GAN/rAC-CT-GAN and cSN-GAN/rcSN-GAN. We resized images from $256 \times 256$ to $64 \times 64$ to shorten the training time.

**Results.** We list the scores in Table 5.[5] The comparison of FID values indicates that the scores depend on the number of data (**noisy**, **mixed** > **clean**) rather than the difference between the baseline and proposed models. This suggests that, in this type noise setting, the scale of the dataset should be made large, even though labels are noisy, to capture an overall distribution. In contrast, the comparison of the GAN-train between the clean and noisy data settings indicates the importance of label accuracy. In the noisy data setting, the scores improve using rAC-GAN or rcGAN but they are still worse than those using AC-GAN and cGAN in the clean data setting. The balanced models are rAC-GAN and rcGAN in the mixed data setting. They are comparable to the models in the noisy data setting in terms of the FID and outperform the models in the clean data setting in terms of the GAN-train. Recently, data augmentation [14, 88] has been studied intensively as an application of conditional generative models. We expect the above findings to provide an important direction in this space.

## 8. Conclusion

Recently, conditional extensions of GANs have shown promise in image generation; however, the limitation here is that they need large-scale accurate class-labeled data to be available. To remedy this, we developed a new family of GANs called rGANs that incorporate a noise transition model into conditional extensions of GANs. In particular, we introduced two variants: rAC-GAN, which is a bridging model between GANs and the noise-robust classification models, and rcGAN, which is an extension of cGAN and solves this problem with no reliance on any classifier. In addition to providing a theoretical background, we demonstrate the effectiveness and limitations of the proposed models through extensive experiments in various settings. In the future, we hope that our findings facilitate the construction of a conditional generative model in real-world scenarios in which only noisy labeled data are available.

---

[4]We tested 10 conditions in total. For each condition, we trained three models with different initializations and report the results averaged over them.

[5]We did not use Intra FID because the number of clean labeled data for each class is few. We did not use the GAN-test because this dataset is challenging and a trained classifier tends to be deceived by noisy data.

## Acknowledgement

## References

[1] D. Angluin and P. Laird. Learning from noisy examples. *Machine Learning*, 2(4):343–370, 1988. 3

[2] P. S. Aritra Ghosh, Himanshu Kumar. Robust loss functions under label noise for deep neural networks. In *AAAI*, 2017. 3

[3] M. Arjovsky and L. Bottou. Towards principled methods for training generative adversarial networks. In *ICLR*, 2017. 2

[4] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In *ICML*, 2017. 2, 27

[5] D. Arpit, S. Jastrzębski, N. Ballas, D. Krueger, E. Bengio, M. S. Kanwal, T. Maharaj, A. Fischer, A. Courville, Y. Bengio, and S. Lacoste-Julien. A closer look at memorization in deep networks. In *ICML*, 2017. 2, 4, 5, 7, 30

[6] A. Bora, E. Price, and A. G. Dimakis. AmbientGAN: Generative models from lossy measurements. In *ICLR*, 2018. 3

[7] A. Brock, J. Donahue, and K. Simonyan. Large scale GAN training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. 1, 2, 22

[8] A. Brock, T. Lim, J. M. Ritchie, and N. Weston. Neural photo editing with introspective adversarial networks. In *ICLR*, 2017. 1

[9] X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, and P. Abbeel. InfoGAN: Interpretable representation learning by information maximizing generative adversarial nets. In *NIPS*, 2016. 5

[10] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo. StarGAN: Unified generative adversarial networks for multi-domain image-to-image translation. In *CVPR*, 2018. 1, 2

[11] H. de Vries, F. Strub, J. Mary, H. Larochelle, O. Pietquin, and A. Courville. Modulating early visual processing by language. In *NIPS*, 2017. 28, 31

[12] E. Denton, S. Chintala, A. Szlam, and R. Fergus. Deep generative image models using a Laplacian pyramid of adversarial networks. In *NIPS*, 2015. 2

[13] V. Dumoulin, J. Shlens, and M. Kudlur. A learned representation for artistic style. In *ICLR*, 2017. 28, 31

[14] M. Frid-Adar, E. Klang, M. Amitai, J. Goldberger, and H. Greenspan. Synthetic data augmentation using GAN for improved liver lesion classification. In *ISBI*, 2018. 1, 8

[15] J. Goldberger and E. Ben-Reuven. Training deep neural-networks using a noise adaptation layer. In *ICLR*, 2017. 3, 4, 5

[16] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014. 1, 2, 3, 4, 5, 27

[17] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville. Improved training of Wasserstein GANs. In *NIPS*, 2017. 2, 3, 6, 22, 27, 28

[18] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On calibration of modern neural networks. In *ICML*, 2017. 30

[19] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. W. Tsang, and M. Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. *arXiv preprint arXiv:1804.06872*, 2018. 3

[20] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 27

[21] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *ECCV*, 2016. 29

[22] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, G. Klambauer, and S. Hochreiter. GANs trained by a two time-scale update rule converge to a Nash equilibrium. In *NIPS*, 2017. 6, 29

[23] S. Iizuka, E. Simo-Serra, and H. Ishikawa. Globally and locally consistent image completion. *ACM Trans. on Graph.*, 36(4):107:1–107:14, 2017. 1

[24] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. 2015. 27

[25] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 1, 2

[26] L. Jiang, Z. Zhou, T. Leung, L.-J. Li, and L. Fei-Fei. MentorNet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *ICML*, 2018. 3

[27] I. Jindal, M. Nokleby, and X. Chen. Learning deep networks from noisy labels with dropout regularization. In *ICDM*, 2016. 3, 4, 5

[28] T. Kaneko, K. Hiramatsu, and K. Kashino. Generative attribute controller with conditional filtered generative adversarial networks. In *CVPR*, 2017. 1, 2

[29] T. Kaneko, K. Hiramatsu, and K. Kashino. Generative adversarial image synthesis with decision tree latent controller. In *CVPR*, 2018. 1, 2

[30] T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. In *ICLR*, 2018. 1, 2

[31] T. Kim, M. Cha, H. Kim, J. K. Lee, and J. Kim. Learning to discover cross-domain relations with generative adversarial networks. In *ICML*, 2017. 1

[32] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 27

[33] D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling. Semi-supervised learning with deep generative models. In *NIPS*, 2014. 2

[34] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 2

[35] A. Krizhevsky. Learning multiple layers of features from tiny images. *Technical report*, 2009. 5

[36] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, 2017. 1, 2

[37] J. Lehtinen, J. Munkberg, J. Hasselgren, S. Laine, T. Karras, M. Aittala, and T. Aila. Noise2Noise: Learning image restoration without clean data. In *ICML*, 2018. 3

[38] J. H. Lim and J. C. Ye. Geometric GAN. *arXiv preprint arXiv:1705.02894*, 2017. 29

[39] M.-Y. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. In *NIPS*, 2017. 1

[40] M. Lucic, K. Kurach, M. Michalski, S. Gelly, and O. Bousquet. Are GANs created equal? A large-scale study. *arXiv preprint arXiv:1711.10337*, 2017. 6, 29

[41] A. Maas, A. Y. Hannun, and A. Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In *ICML Workshop*, 2013. 27

[42] E. Malach and S. Shalev-Shwartz. Decoupling "when to update" from "how to update". In *NIPS*, 2017. 3

[43] E. Mansimov, E. Parisotto, J. L. Ba, and R. Salakhutdinov. Generating images from captions with attention. In *ICLR*, 2016. 2

[44] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. P. Smolley. Least squares generative adversarial networks. In *ICCV*, 2017. 2

[45] A. Menon, B. van Rooyen, C. S. Ong, and B. Williamson. Learning from corrupted binary labels via class-probability estimation. In *ICML*, 2015. 30

[46] L. Mescheder, A. Geiger, and S. Nowozin. Which training methods for GANs do actually converge? In *ICML*, 2018. 2

[47] M. Mirza and S. Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 1, 2, 3, 4, 6, 27

[48] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral normalization for generative adversarial networks. In *ICLR*, 2018. 1, 2, 6, 22, 27, 28, 30, 31

[49] T. Miyato and M. Koyama. cGANs with projection discriminator. In *ICLR*, 2018. 1, 2, 3, 4, 6, 28, 29

[50] V. Nair and G. E. Hinton. Rectified linear units improve restricted Boltzmann machines. In *ICML*, 2010. 27

[51] N. Natarajan, I. S. Dhillon, P. K. Ravikumar, and A. Tewari. Learning with noisy labels. In *NIPS*, 2013. 3

[52] A. Odena, C. Olah, and J. Shlens. Conditional image synthesis with auxiliary classifier GANs. In *ICML*, 2017. 1, 2, 3, 27

[53] G. Patrini, A. Rozza, A. K. Menon, R. Nock, and L. Qu. Making deep neural networks robust to label noise: A loss correction approach. In *CVPR*, 2017. 3, 4, 5, 6, 7, 8, 18, 30

[54] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2016. 2, 6, 22, 27

[55] S. Reed, Z. Akata, S. Mohan, S. Tenka, B. Schiele, and H. Lee. Learning what and where to draw. In *NIPS*, 2016. 2

[56] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee. Generative adversarial text to image synthesis. In *ICML*, 2016. 2

[57] S. Reed, H. Lee, D. Anguelov, C. Szegedy, D. Erhan, and A. Rabinovich. Training deep neural networks on noisy labels with bootstrapping. In *ICLR*, 2015. 3

[58] S. Reed, A. van den Oord, N. Kalchbrenner, V. Bapst, M. Botvinick, and N. de Freitas. Generating interpretable

[59] M. D. Reid and R. C. Williamson. Composite binary losses. *JMLR*, 11:2387–2422, 2010. 4

[60] M. Ren, W. Zeng, B. Yang, and R. Urtasun. Learning to reweight examples for robust deep learning. In *ICML*, 2018. 3

[61] D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *ICML*, 2014. 2

[62] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training GANs. In *NIPS*, 2016. 2, 29

[63] K. Shmelkov, C. Schmid, and K. Alahari. How good is my GAN? In *ECCV*, 2018. 6, 29, 30

[64] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb. Learning from simulated and unsupervised images through adversarial training. In *CVPR*, 2017. 1

[65] S. Sukhbaatar, J. Bruna, M. Paluri, L. Bourdev, and R. Fergus. Training convolutional networks with noisy labels. In *ICLR Workshop*, 2015. 3, 4, 5

[66] D. Tanaka, D. Ikami, T. Yamasaki, and K. Aizawa. Joint optimization framework for learning with noisy labels. In *CVPR*, 2018. 3, 4, 22

[67] L. Theis, A. van den Oord, and M. Bethge. A note on the evaluation of generative models. In *ICLR*, 2016. 6

[68] D. Tran, R. Ranganath, and D. M. Blei. Deep and hierarchical implicit models. *arXiv preprint arXiv:1702.08896*, 2017. 29

[69] A. van den Oord, N. Kalchbrenner, and K. Kavukcuoglu. Pixel recurrent neural networks. In *ICML*, 2016. 2

[70] A. van den Oord, N. Kalchbrenner, O. Vinyals, L. Espeholt, A. Graves, and K. Kavukcuoglu. Conditional image generation with pixelCNN decoders. In *NIPS*, 2016. 2

[71] B. van Rooyen, A. Menon, and R. C. Williamson. Learning with symmetric label noise: The importance of being unhinged. In *NIPS*, 2015. 6

[72] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, G. Liu, A. Tao, J. Kautz, and B. Catanzaro. Video-to-video synthesis. *arXiv preprint arXiv:1808.06601*, 2018. 1, 2

[73] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro. High-resolution image synthesis and semantic manipulation with conditional GANs. In *CVPR*, 2018. 1, 2

[74] X. Wei, B. Gong, Z. Liu, W. Lu, and L. Wang. Improving the improved training of Wasserstein GANs: A consistency term and its dual effect. In *ICLR*, 2018. 2, 6, 27, 28

[75] T. Xiao, T. Xia, Y. Yang, C. Huang, and X. Wang. Learning from massive noisy labeled data for image classification. In *CVPR*, 2015. 2, 8, 21

[76] B. Xu, N. Wang, T. Chen, and M. Li. Empirical evaluation of rectified activations in convolutional network. In *ICML Workshop*, 2015. 27

[77] T. Xu, P. Zhang, Q. Huang, H. Zhang, Z. G. X. Huang, and X. He. AttnGAN: Fine-grained text to image generation with attentional generative adversarial networks. In *CVPR*, 2018. 2

[78] X. Yan, J. Yang, K. Sohn, and H. Lee. Attribute2image: Conditional image generation from visual attributes. In *ECCV*, 2016. 2

[79] Z. Yi, H. Zhang, P. Tan, and M. Gong. DualGAN: Unsupervised dual learning for image-to-image translation. In *ICCV*, 2017. 1

[80] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning requires rethinking generalization. In *ICLR*, 2017. 2, 4, 5, 30

[81] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz. mixup: Beyond empirical risk minimization. In *ICLR*, 2018. 4, 29

[82] H. Zhang, I. Goodfellow, D. Metaxas, and A. Odena. Self-attention generative adversarial networks. *arXiv preprint arXiv:1805.08318*, 2018. 1, 2

[83] H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, and D. Metaxas. StackGAN++: Realistic image synthesis with stacked generative adversarial networks. *arXiv preprint arXiv:1710.10916*, 2017. 2

[84] H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, and D. N. Metaxas. StackGAN: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *ICCV*, 2017. 2

[85] Z. Zhang and M. R. Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. *arXiv preprint arXiv:1805.07836*, 2018. 3

[86] Z. Zhang, Y. Song, and H. Qi. Age progression/regression by conditional adversarial autoencoder. In *CVPR*, 2017. 1, 2

[87] J. Zhao, M. Mathieu, and Y. LeCun. Energy-based generative adversarial network. In *ICLR*, 2017. 2

[88] Z. Zhong, L. Zheng, Z. Zheng, S. Li, and Y. Yang. Camera style adaptation for person re-identification. In *CVPR*, 2018. 1, 8

[89] J.-Y. Zhu, P. Krähenbühl, E. Shechtman, and A. A. Efros. Generative visual manipulation on the natural image manifold. In *ECCV*, 2016. 1

[90] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 1

# A. Contents

# B. Extended results

## B.1. Extended results of Section 7.1

As extended results of Section 7.1 (comprehensive study), we provide the image samples generated using the models evaluated in Section 7.1 in Figures 6–9. An outline of the content is as follows:

- Figure 6: Image samples generated using DCGANs (AC-DCGAN, rAC-DCGAN, cDCGAN, and rcDC-GAN) on CIFAR-10
- Figure 7: Image samples generated using WGAN-GPs (AC-WGAN-GP, rAC-WGAN-GP, cWGAN-GP, and rcWGAN-GP) on CIFAR-10
- Figure 8: Image samples generated using CT-GANs (AC-CT-GAN, rAC-CT-GAN, cCT-GAN, and rcCT-GAN) on CIFAR-10
- Figure 9: Image samples generated using SN-GANs (AC-SN-GAN, rAC-SN-GAN, cSN-GAN, and rcSN-GAN) on CIFAR-10

## B.2. Extended results of Section 7.2

We provide the extended results of Section 7.2 (effects of estimated $T'$) in Table 6, Figure 10, and Figure 11. An outline of the content is as follows:

- Table 6: Quantitative results using estimated $T'$ (extended version of Table 3)
- Figure 10: Visualization of Table 6
- Figure 11: Image samples generated using rAC-CT-GAN with estimated $T'$ and rcSN-GAN with estimated $T'$ on CIFAR-10

## B.3. Extended results of Section 7.3

We provide the extended results of Section 7.3 (evaluation of the improved technique) in Table 7 and Figure 12. An outline of the content is as follows:

- Table 7: Quantitative results using the improved technique (extended version of Table 4)
- Figure 12: Image samples generated using the improved rAC-GANs and improved rcGANs with combinations of DCGAN, WGAN-GP, CT-GAN, and SN-GAN on CIFAR-10

## B.4. Extended results of Section 7.4

We provide the extended results of Section 7.4 (evaluation on the real-world noise) in Figure 13. An outline of the content is as follows:

- Figure 13: Image samples generated using AC-CT-GAN, rAC-CT-GAN, cSN-GAN, and rcSN-GAN on Clothing1M (clean, noisy, and mixed settings)
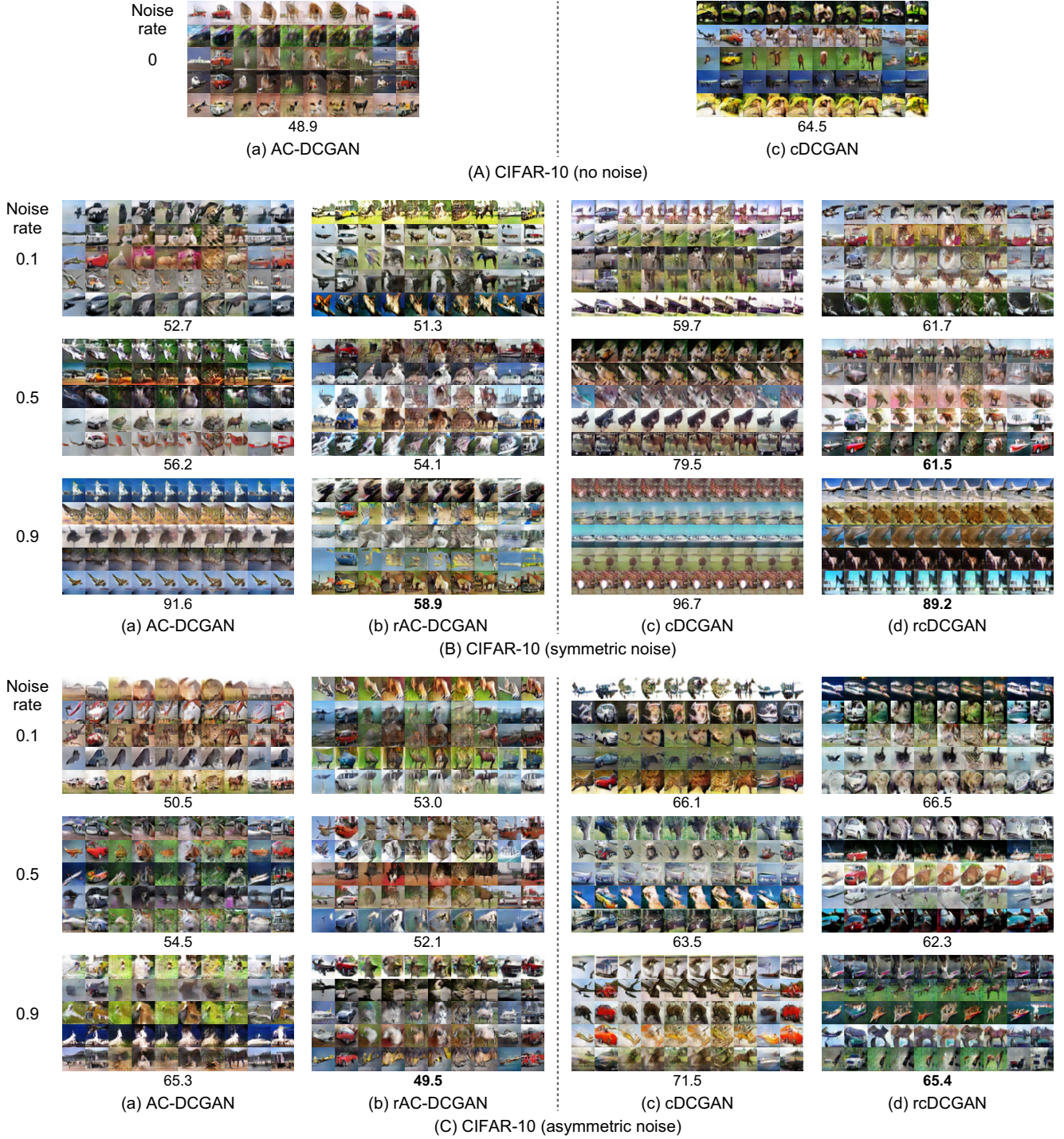
Figure 6. Image samples generated using (a) AC-DCGAN, (b) rAC-DCGAN, (c) cDCGAN, and (d) rcDCGAN on CIFAR-10 ((A) no noise, (B) symmetric noise, and (C) asymmetric noise). These models are discussed in Section 7.1. In each picture block, each column shows samples associated with the same class. Each row includes samples generated from a fixed $z$ and a varied $y^g$. The value below each picture block represents the achieved Intra FID (which is the same as the value reported in Figure 4). The smaller the value, the better. When the score difference between the baseline models (AC-DCGAN and cDCGAN) and the proposed models (rAC-DCGAN and rcDCGAN) is more than 3 points, we use bold font to indicate the better model.

Figure 7. Image samples generated using (a) AC-WGAN-GP, (b) rAC-WGAN-GP, (c) cWGAN-GP, and (d) rcWGAN-GP on CIFAR-10 ((A) no noise, (B) symmetric noise, and (C) asymmetric noise). These models are discussed in Section 7.1. In each picture block, each column shows samples associated with the same class. Each row includes samples generated from a fixed $z$ and a varied $y^g$. The value below each picture block represents the achieved Intra FID (which is the same as the value reported in Figure 4). The smaller the value, the better. When the score difference between the baseline models (AC-WGAN-GP and cWGAN-GP) and the proposed models (rAC-WGAN-GP and rcWGAN-GP) is more than 3 points, we use bold font to indicate the better model.

14

Figure 8. Image samples generated using (a) AC-CT-GAN, (b) rAC-CT-GAN, (c) cCT-GAN, and (d) rcCT-GAN on CIFAR-10 ((A) no noise, (B) symmetric noise, and (C) asymmetric noise). These models are discussed in Section 7.1. In each picture block, each column shows samples associated with the same class. Each row includes samples generated from a fixed $z$ and a varied $y^g$. The value below each picture block represents the achieved Intra FID (which is the same as the value reported in Figure 4). The smaller the value, the better. When the score difference between the baseline models (AC-CT-GAN and cCT-GAN) and the proposed models (rAC-CT-GAN and rcCT-GAN) is more than 3 points, we use bold font to indicate the better model.
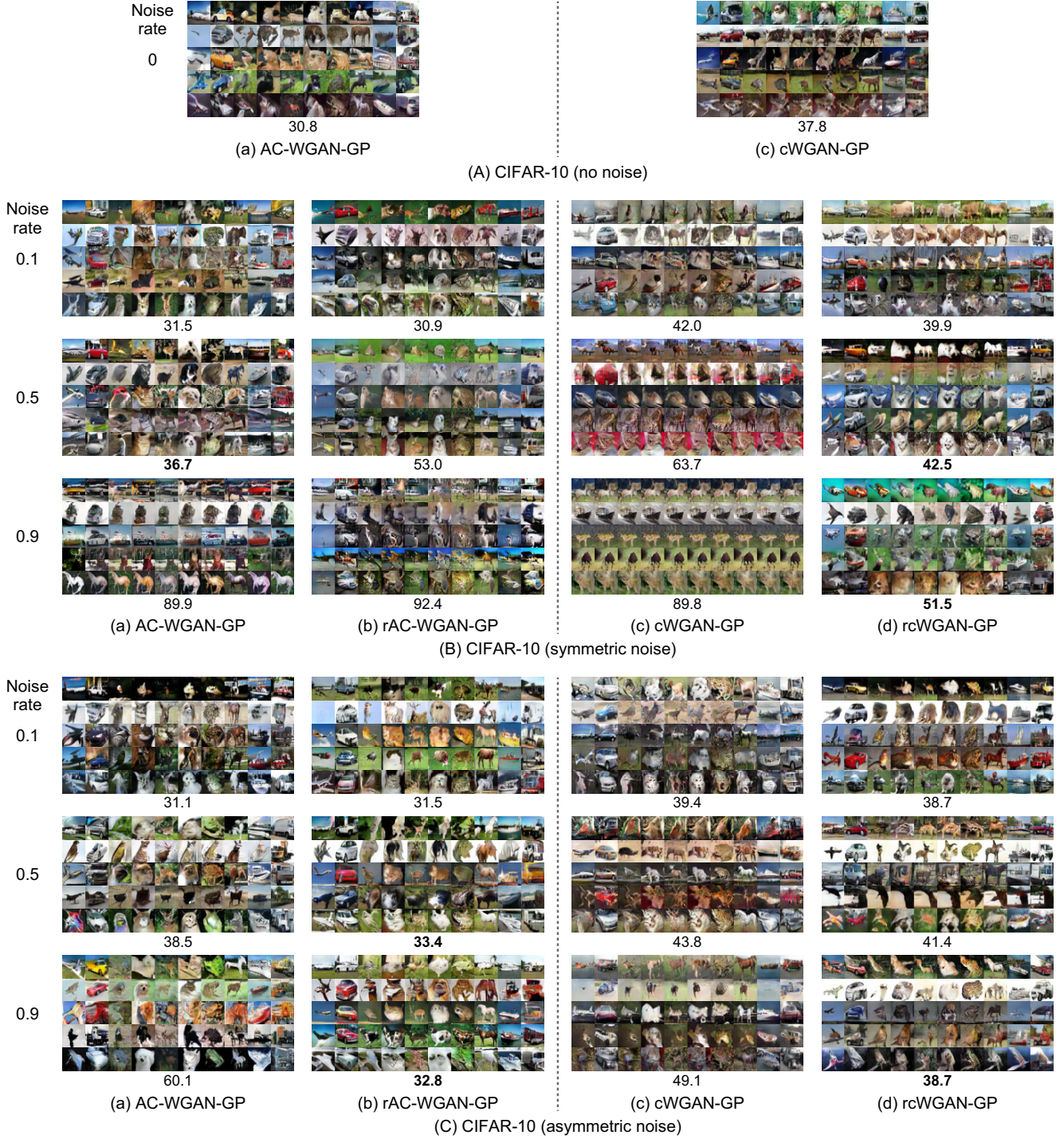
15

Figure 9. Image samples generated using (a) AC-SN-GAN, (b) rAC-SN-GAN, (c) cSN-GAN, and (d) rcSN-GAN on CIFAR-10 ((A) no noise, (B) symmetric noise, and (C) asymmetric noise). These models are discussed in Section 7.1. In each picture block, each column shows samples associated with the same class. Each row includes samples generated from a fixed $z$ and a varied $y^g$. The value below each picture block represents the achieved Intra FID (which is the same as the value reported in Figure 4). The smaller the value, the better. When the score difference between the baseline models (AC-SN-GAN and cSN-GAN) and the proposed models (rAC-SN-GAN and rcSN-GAN) is more than 3 points, we use bold font to indicate the better model.
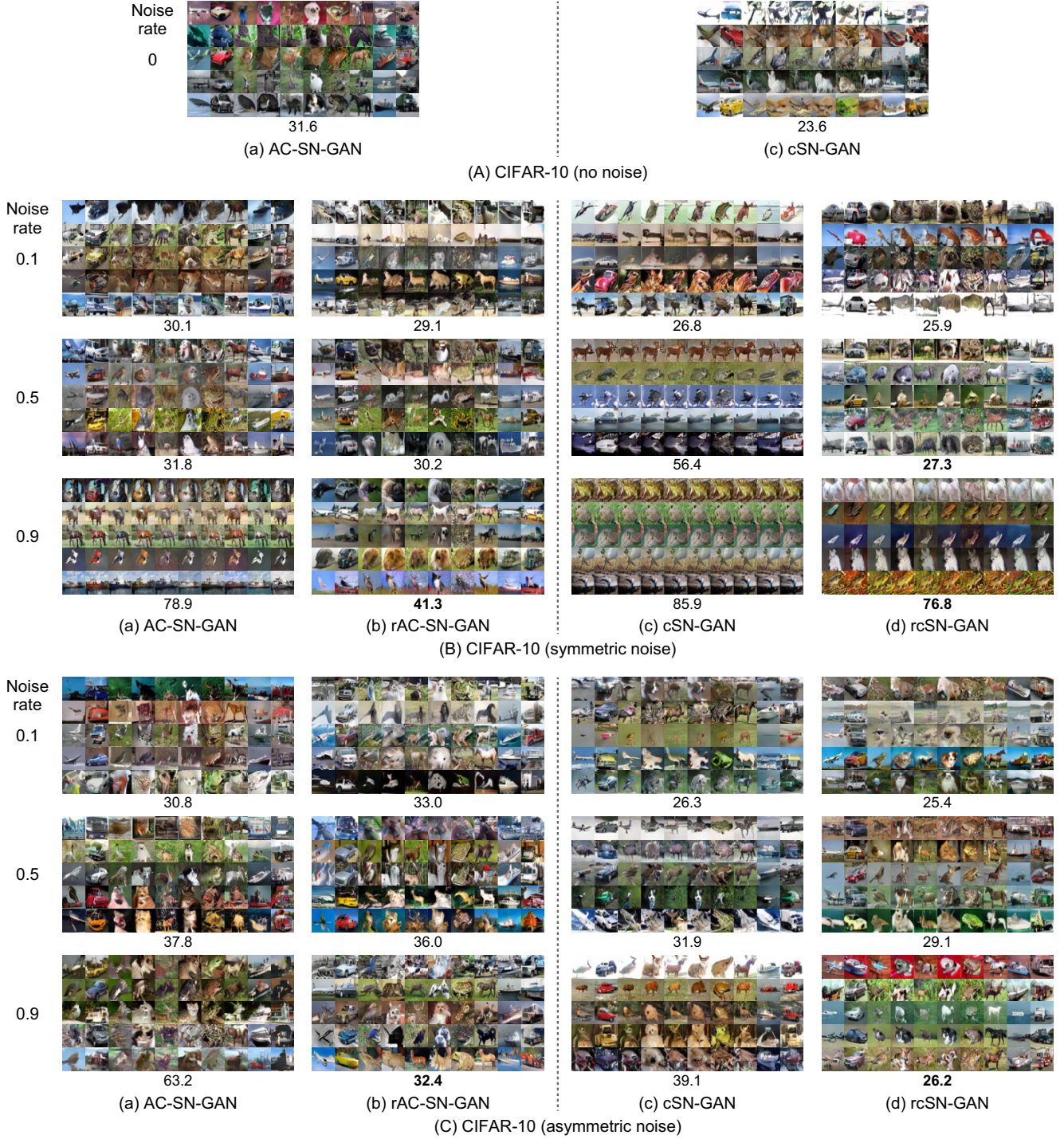
**Table (a) CIFAR-10**

| Model | Metric | CIFAR-10 (symmetric noise) | | | | | CIFAR-10 (asymmetric noise) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 0.1 | 0.3 | 0.5 | 0.7 | 0.9 | 0.1 | 0.3 | 0.5 | 0.7 | 0.9 |
| rAC-CT-GAN with $T'$ (AC-CT-GAN) | FID ↓ | 10.9 | 11.4 | 11.3 | 11.5 | 13.0 | 10.8 | 10.2 | 10.2 | 10.4 | 11.0 |
| | | (11.4) | (13.4) | (14.0) | (13.5) | (14.3) | (10.8) | (11.0) | (11.0) | (11.0) | (10.9) |
| | Intra FID ↓ | 28.7 | **31.0** | **30.1** | **31.7** | **38.9** | 28.5 | **27.4** | **31.2** | **35.0** | **36.8** |
| | | (29.8) | (35.1) | (37.4) | (36.4) | (61.9) | (28.3) | (30.7) | (35.4) | (45.7) | (62.6) |
| | GAN-test ↑ | 95.3 | 93.2 | **92.0** | 87.7 | **70.4** | 94.9 | 92.9 | **85.2** | **78.5** | **76.6** |
| | | (94.7) | (91.7) | (88.9) | (86.7) | (40.9) | (94.0) | (91.0) | (78.8) | (69.0) | (62.7) |
| | GAN-train ↑ | 78.7 | **75.9** | **76.9** | **73.7** | **63.4** | 79.8 | **79.5** | **74.0** | **69.1** | **67.3** |
| | | (78.1) | (72.0) | (70.7) | (67.9) | (34.5) | (78.7) | (74.1) | (62.5) | (51.5) | (47.7) |
| rcSN-GAN with $T'$ (cSN-GAN) | FID ↓ | 10.7 | 11.9 | 12.4 | 12.1 | 15.0 | 10.8 | 10.8 | 11.0 | 10.9 | 11.3 |
| | | (11.0) | (12.9) | (14.7) | (14.8) | (14.8) | (11.2) | (10.9) | (11.4) | (10.7) | (11.0) |
| | Intra FID ↓ | 25.5 | **29.4** | **29.4** | **29.7** | 87.4 | 25.7 | 26.0 | **28.7** | 32.6 | **33.9** |
| | | (26.8) | (38.5) | (56.4) | (72.0) | (85.9) | (26.3) | (28.2) | (31.9) | (33.7) | (39.1) |
| | GAN-test ↑ | **85.3** | **79.0** | **84.8** | **82.8** | 15.9 | 86.6 | **87.2** | **84.0** | **74.9** | **71.2** |
| | | (81.1) | (60.2) | (38.5) | (23.2) | (13.1) | (85.1) | (77.3) | (70.8) | (66.3) | (59.5) |
| | GAN-train ↑ | 80.7 | **78.1** | **77.4** | **75.6** | 15.0 | 80.5 | **79.0** | **75.7** | **69.3** | **65.7** |
| | | (79.5) | (69.2) | (45.5) | (28.5) | (14.5) | (80.4) | (73.5) | (65.9) | (59.8) | (51.0) |

(a) CIFAR-10

**Table (b) CIFAR-100**

| Model | Metric | CIFAR-100 (symmetric noise) | | | | | CIFAR-100 (asymmetric noise) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 0.1 | 0.3 | 0.5 | 0.7 | 0.9 | 0.1 | 0.3 | 0.5 | 0.7 | 0.9 |
| rAC-CT-GAN with $T'$ (AC-CT-GAN) | FID ↓ | 19.7 | 19.3 | 17.7 | 17.3 | 18.5 | 19.4 | 19.3 | 19.7 | 18.8 | 19.0 |
| | | (19.2) | (19.1) | (18.7) | (18.0) | (18.0) | (18.9) | (18.5) | (19.2) | (19.6) | (19.3) |
| | GAN-test ↑ | **76.6** | 67.1 | **68.1** | *1.0* | 2.5 | 74.1 | 68.9 | 28.7 | 7.2 | 2.2 |
| | | (72.4) | (65.0) | (63.1) | (48.0) | (9.1) | (75.5) | (68.4) | (34.4) | (8.7) | (3.8) |
| | GAN-train ↑ | 21.2 | 21.4 | 23.3 | *1.0* | 2.3 | 19.1 | 19.9 | 10.7 | 5.5 | 3.9 |
| | | (21.7) | (22.8) | (21.7) | (19.3) | (5.1) | (21.4) | (20.8) | (12.2) | (5.8) | (4.0) |
| rcSN-GAN with $T'$ (rcSN-GAN) | FID ↓ | 14.3 | 16.6 | 17.5 | 20.0 | 19.8 | 13.8 | 14.1 | 14.7 | 14.7 | 13.9 |
| | | (14.2) | (16.9) | (18.9) | (19.4) | (18.7) | (13.3) | (14.2) | (14.6) | (14.4) | (13.5) |
| | GAN-test ↑ | 53.4 | 36.6 | **37.7** | *1.0* | 1.7 | **65.0** | **63.0** | 32.4 | *7.8* | 3.8 |
| | | (54.3) | (33.9) | (13.9) | (5.9) | (1.9) | (56.1) | (41.8) | (27.5) | (15.6) | (5.4) |
| | GAN-train ↑ | 40.1 | 32.8 | **31.3** | *1.0* | 1.8 | 41.7 | **39.3** | 20.1 | *6.1* | 3.9 |
| | | (39.7) | (33.2) | (16.9) | (7.7) | (1.9) | (41.7) | (33.3) | (20.7) | (11.1) | (4.8) |

(b) CIFAR-100

Table 6. Extended version of Table 3. Quantitative results using the estimated $T'$. These results are discussed in Section 7.2. In each table, the second row indicates a noise rate $\mu \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$. Under the third row, each odd row contains the scores for the proposed models (i.e., rAC-CT-GAN or rcSN-GAN) with $T'$ and each even row (denoted in parenthesis) includes the scores for the baseline models (i.e., AC-CT-GAN or cSN-GAN). Bold and italic fonts indicate that the score for the proposed models is better or worse by more than 3 points than that for the baseline models, respectively. See also Figure 10 that visualizes this information as graphs.
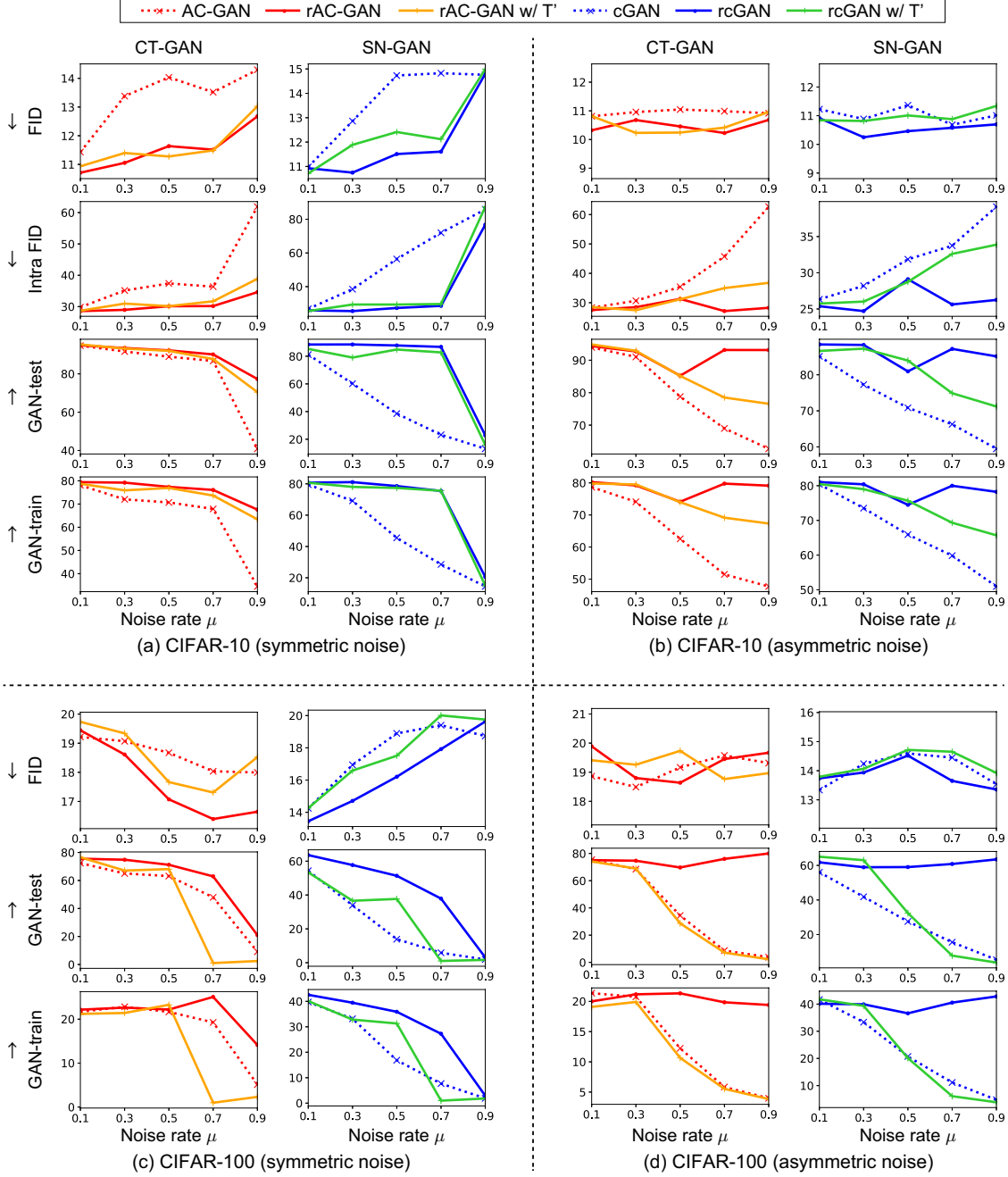
Figure 10. Visualization of Table 6. Comparison of the quantitative results using the baseline models (AC-CT-GAN and cSN-GAN), the proposed models (rAC-CT-GAN and rcSN-GAN) with the known $T$, and the proposed models with the estimated $T'$. The scale is adjusted on each graph for easy viewing. As discussed in Section 7.2, in CIFAR-10, even using $T'$, rAC-CT-GAN and rcSN-GAN outperform conventional AC-CT-GAN and cSN-GAN, respectively, and show robustness to label noise. Furthermore, rAC-CT-GAN and rcSN-GAN with $T$ and those with $T'$ are almost similar except for the asymmetric noise with a higher noise rate (i.e., 0.7 and 0.9). In CIFAR-100, when the noise rate is low, rAC-CT-GAN and rcSN-GAN work moderately well; however, in highly noisy settings, their performance is degraded. This implies the limitation of estimating $T'$ from the data in which there is a high-rate mixture and there is a limited number of images per class (500). This is also mentioned in the previous study [53]. Further improvement remains as an open issue. The precise values for this figure are provided in Table 6.
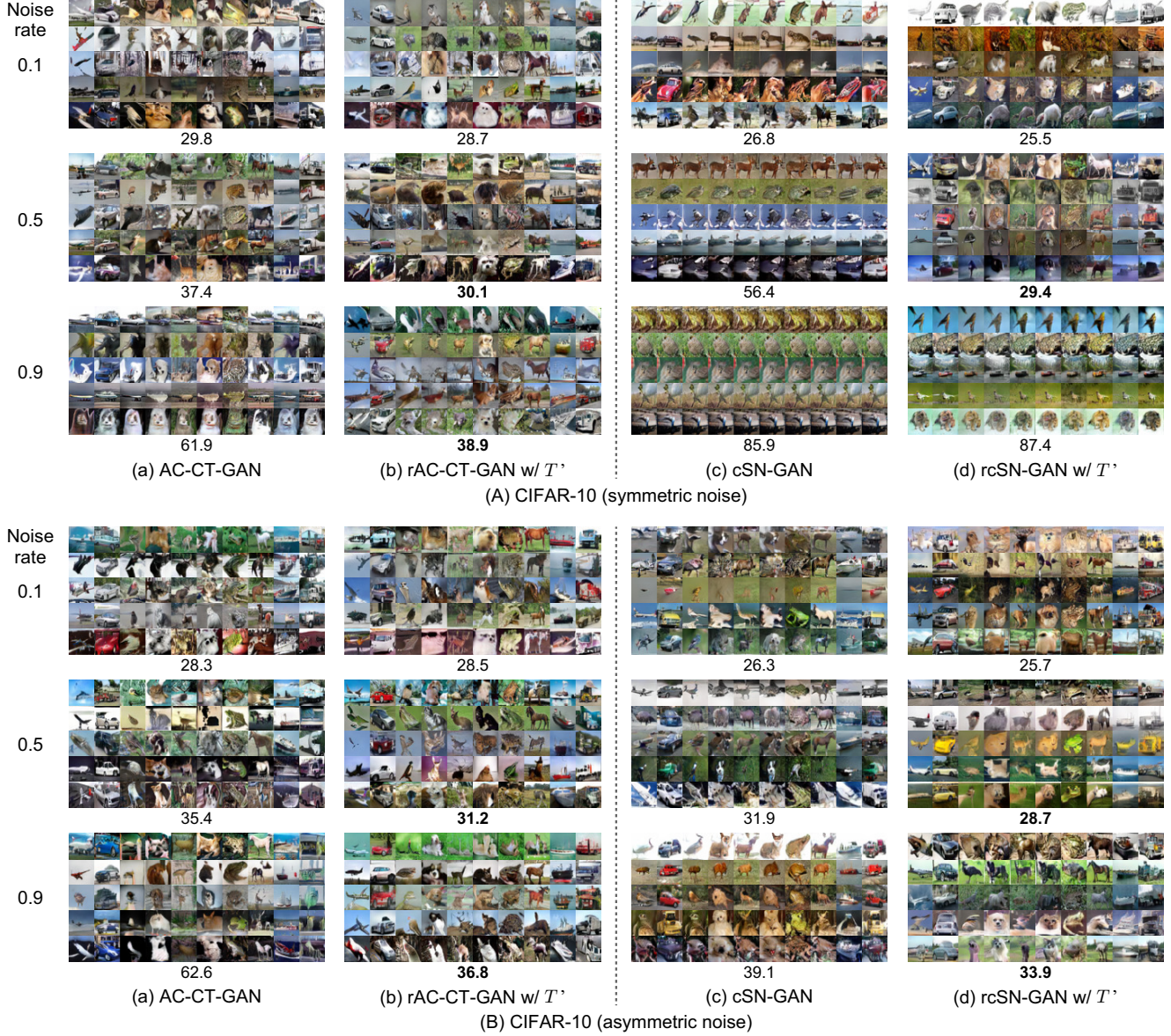
Noise rate

0.1

29.8 28.7 26.8 25.5

0.5

37.4 **30.1** 56.4 **29.4**

0.9

61.9 **38.9** 85.9 87.4

(a) AC-CT-GAN  (b) rAC-CT-GAN w/ $T'$  (c) cSN-GAN  (d) rcSN-GAN w/ $T'$

(A) CIFAR-10 (symmetric noise)

Noise rate

0.1

28.3 28.5 26.3 25.7

0.5

35.4 **31.2** 31.9 **28.7**

0.9

62.6 **36.8** 39.1 **33.9**

(a) AC-CT-GAN  (b) rAC-CT-GAN w/ $T'$  (c) cSN-GAN  (d) rcSN-GAN w/ $T'$

(B) CIFAR-10 (asymmetric noise)

Figure 11. Image samples generated using (a) AC-CT-GAN, (b) rAC-CT-GAN with the estimated $T'$, (c) cSN-GAN, and (d) rcSN-GAN with the estimated $T'$ on CIFAR-10 ((A) symmetric noise and (B) asymmetric noise). These models are discussed in Section 7.2. In each picture block, each column shows samples associated with the same class. Each row includes samples generated from a fixed $z$ and a varied $y^g$. The value below each picture block represents the achieved Intra FID (which is the same as the value reported in Tables 3 and 6). When the score difference between the baseline models ((a) AC-CT-GAN and (c) cSN-GAN) and the proposed models ((b) rAC-CT-GAN with $T'$ and (d) rcSN-GAN with $T'$) is more than 3 points, we use bold font to indicate the better model. Refer to Figures 8 and 9 for comparison with rAC-CT-GAN with the known $T$ and rcSN-GAN with the known $T$, respectively.

19

| Model | Metric | CIFAR-10 (symmetric noise) | | | | CIFAR-100 (symmetric noise) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | A | B | C | D | A | B | C | D |
| Improved rAC-GAN (rAC-GAN) | FID ↓ | 27.9 (28.7) | **14.7** (38.8) | 12.4 (12.7) | 13.5 (13.5) | 33.1 (34.2) | **20.4** (41.1) | 17.2 (16.6) | 18.4 (18.6) |
| | Intra FID ↓ | **55.7** (58.9) | **34.6** (92.4) | 33.4 (34.6) | **36.9** (41.3) | – – | – – | – – | – – |
| | GAN-test ↑ | 65.1 (62.7) | **77.7** (27.1) | 78.2 (77.3) | **63.5** (59.6) | 26.2 (27.2) | **22.5** (1.0) | 21.5 (21.0) | **15.4** (7.9) |
| | GAN-train ↑ | 59.9 (58.7) | **70.8** (26.3) | 69.1 (67.6) | **59.7** (51.9) | 17.1 (17.4) | **16.3** (1.0) | 14.8 (14.1) | **11.7** (6.9) |
| Improved rcGAN (rcGAN) | FID ↓ | **30.4** (35.4) | **16.9** (22.2) | 14.2 (14.8) | 14.9 (14.8) | *50.2* (40.1) | **25.8** (31.2) | 18.0 (17.5) | 18.7 (19.6) |
| | Intra FID ↓ | **76.9** (89.2) | **39.6** (51.5) | **52.9** (63.2) | **48.2** (76.8) | – – | – – | – – | – – |
| | GAN-test ↑ | 27.3 (18.9) | 65.7 (58.4) | 38.9 (31.2) | 48.8 (22.8) | 4.5 (1.5) | 12.0 (9.2) | 9.5 (5.3) | 6.1 (3.1) |
| | GAN-train ↑ | 31.9 (25.3) | 60.7 (48.7) | 36.7 (30.6) | 47.3 (20.5) | 6.0 (1.7) | 10.3 (8.3) | 7.5 (4.7) | 4.4 (2.9) |

Table 7. Extended version of Table 4. Quantitative results using the improved technique. These results are discussed in Section 7.3. In the second row, A, B, C, and D indicate DCGAN, WGAN-GP, CT-GAN, and SN-GAN, respectively. We evaluated the models in severely noisy settings (i.e., symmetric noise with a noise rate 0.9). Under the third row, each odd row contains the scores for the improved rAC-GAN or improved rcGAN and each even row (denoted in parenthesis) contains the scores for the naive rAC-GAN or naive rcGAN. Bold and italic fonts indicate that the score for the improved models is better or worse by more than 3 points than that for the naive models, respectively.
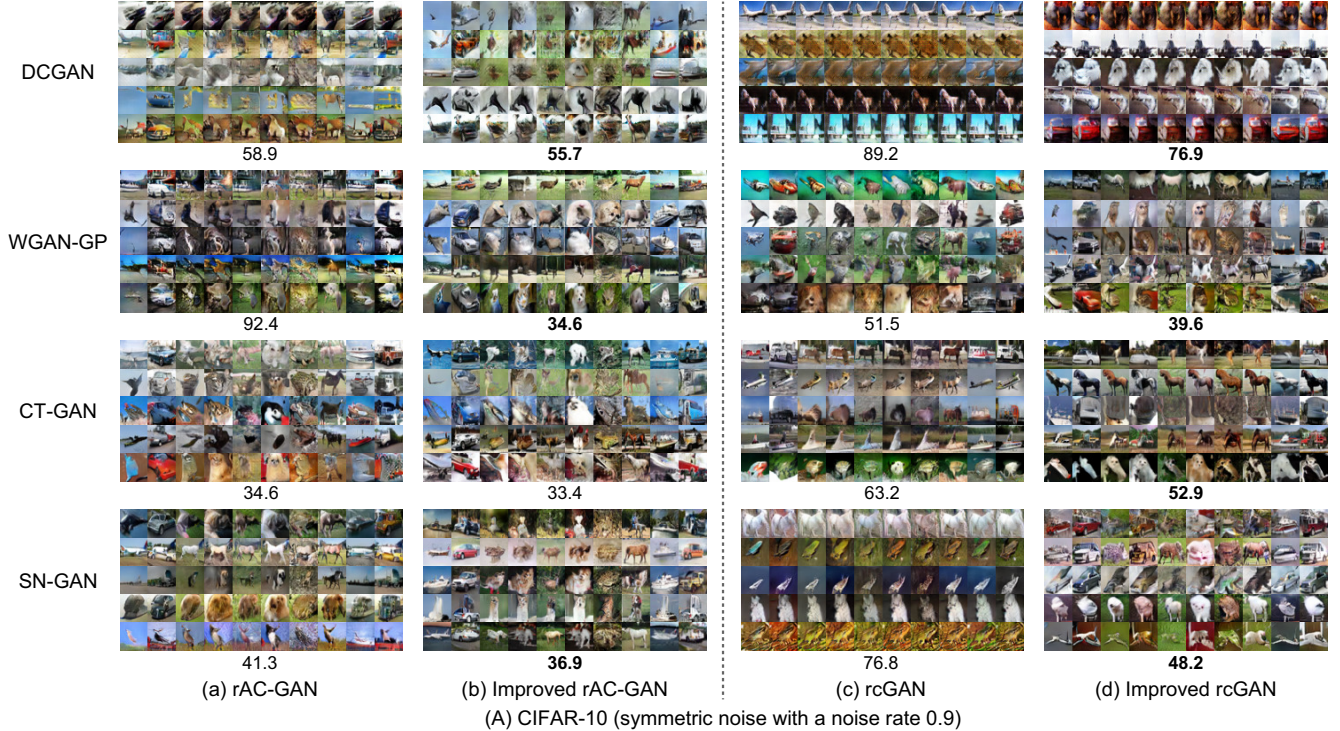


Figure 12. Image samples generated using (a) rAC-GAN, (b) improved rAC-GAN, (c) rcGAN, and (d) improved rcGAN on CIFAR-10 in severely noisy settings (i.e., symmetric noise with a noise rate 0.9). These models are discussed in Section 7.3. In each picture block, each column shows samples associated with the same class. Each row includes samples generated from a fixed $z$ and a varied $y^g$. The value below each picture block represents the achieved Intra FID (which is the same as the value reported in Tables 4 and 7). When the score difference between the naive models ((a) rAC-GAN and (c) rcGAN) and the improved models ((b) improved rAC-GAN and (d) improved rcGAN) is more than 3 points, we use bold font to indicate the better model.

Figure 13. Image samples generated using (a) AC-CT-GAN, (b) rAC-CT-GAN, (c) cSN-GAN, and (d) rcSN-GAN on Clothing1M ((A) clean, (B) noisy, and (C) mixed settings). These models are discussed in Section 7.4. In each picture block, each column shows samples belonging to the same class. From left to right, each column represents t-shirt, shirt, knitwear, chiffon, sweater, hoodie, windbreaker, jacket, down coat, suit, shawl, dress, vest, and underwear, respectively. Each row includes samples generated from a fixed $z$ and a varied $y^g$. The value below each picture block represents the achieved GAN-train (which is the same as the value reported in Table 5). The larger the value, the better. Bold font indicates a better score in each block. Note that this dataset is challenging (annotation accuracy is only 61.54% [75]) and correct labeling is also difficult for humans.

## C. Additional analysis

### C.1. Effect of gap between real and model $T$

In Section 7.2, we evaluated the models with the estimated $T'$ and examined the effect when there is a gap between the real $T$ and model $T$ (particularly, $T'$ in this case). To further investigate such an effect, we conducted an additional experiment. In the following, to clarify the difference, we denote the real $T$ and model $T$ by $T^r$ and $T^g$, respectively, and their corresponding noise rates by $\mu^r$ and $\mu^g$, respectively. In Section 7.1, we examined the performance change when $\mu^r$ and $\mu^g$ are varied at the same time (i.e., $\mu^r = \mu^g$). In contrast, in this section, to inspect the effect of the gap between $T^r$ and $T^g$, we fixed $\mu^r$ as a constant value ($\mu^r = 0$ or $\mu^r = 0.5$) and investigated the performance change when $\mu^g$ is varied ($\mu^g \in \{0, 0.1, 0.3, 0.5, 0.7, 0.9\}$). Figures 14 and 15 show the results for noise rates $\mu^r = 0$ and $\mu^r = 0.5$, respectively. Although there is a dependency on the models, datasets, and evaluation metrics, we find that the quantity degradation is relatively small when the gap between $\mu^r$ and $\mu^g$ is within $\pm 0.2$. However, in this situation, the theoretical guarantees supported by Theorems 1 and 2 do not hold, and we admit that there is room to explore these observations theoretically in future work.

### C.2. Effect of learning rate

Recent studies (e.g., [66]) show that a high learning rate is useful for preventing a classifier DNN from memorizing noisy labels. To explore such an effect on conditional generative models, we performed a comparative study using the models with different learning rates. In particular, we evaluated the baseline models (i.e., AC-CT-GAN and cSN-GAN) and the proposed models (i.e., rAC-CT-GAN and rcSN-GAN) in severely noisy settings (i.e., symmetric noise with a noise rate 0.9). We selected the initial learning rate $\alpha$ from 0.0001, 0.0002, 0.0004, and 0.0008. As described in Appendix D, the default parameter of $\alpha$ is 0.0002. We trained the models for $100k$ generator iterations and decayed $\alpha$ to 0 over $100k$ iterations in all settings.

**Results.** We display the results in Figure 16. As discussed in the previous studies, generally the GAN training itself is not stable and has sensitivity to the learning rate (e.g., the authors of DCGAN [54] recommended a low learning rate). Therefore, the relationship between the model and the learning rate in GANs might be more difficult to explain than that in the classifier DNNs. However, we observed two tendencies through the experiments:

- As the learning rate increases (particularly ranged from 0.0001 to 0.0004), the quantitative scores tend to become better in the proposed models; however, such benefits are small in the baseline models (particularly in cSN-GAN). We argue that this is because our pro-

posed models can employ noisy labels as useful conditional information and this allows for suppressing the training instability resulting from a high learning rate.
- However, when using an extensively high learning rate (e.g., 0.0008), the scores degrade even when using the proposed models (particularly when using rAC-CT-GAN). This implies the necessity of a careful parameter tuning.

As per the latest studies (e.g., [17, 48]), the dependency on hyperparameter settings is being improved, and we expect that a more label-noise robust model will be constructed along with the advances in GANs.

### C.3. Effect of batch size

Another important factor with regard to the training is the batch size. In particular, it might be critical in noisy label settings because, as the batch size becomes small, the factors for distinguishing between right and wrong labels also become fewer. To investigate this effect, we conducted a comparative study using the models with different batch sizes. We selected the batch sizes from 32, 64, and 128. As described in Appendix D, the default batch size is 64. In this analysis, we set the learning rate $\alpha$ to 0.0002 (default).

**Results.** We show the results in Figure 17. As was the case with the learning rate, the batch size affects the GAN training itself. Therefore, it is not easy to explain precisely the relationship between the model and the batch size. However, we observed a similar tendency to that of the learning rate, i.e., the proposed models benefit from an increasing batch size, whereas such benefits are small in the baseline models. The latest study [7] demonstrates that, by incorporating some techniques, it is possible to obtain GAN training stability even when using a large batch size (e.g., a batch size of 2048). We expect that the performance of rAC-GAN and rcGAN will be improved along with such advances.

### C.4. Distance to noisy labeled data

In the main text, we used Intra FID to measure the distance between the generated data distribution and the *clean* labeled data distribution. Another interesting metric is the distance between the generated data distribution and the *noisy* labeled data distribution. To assess it, we computed Intra FID between the samples generated by $G$ and the real samples belonging to the class of concern in terms of *noisy* labels. We show the results in Figure 18.[6] These results support the finding, discussed in the last paragraph in Section 7.1, i.e., cGAN can fit even noisy labels, and AC-GAN shows robustness for symmetric noise. We found that these tendencies occur independently of the GAN configurations.

---

[6]We calculated these scores only for CIFAR-10 with symmetric noise because in the other settings the number of noisy labeled data for each class is insufficient to use this metric.
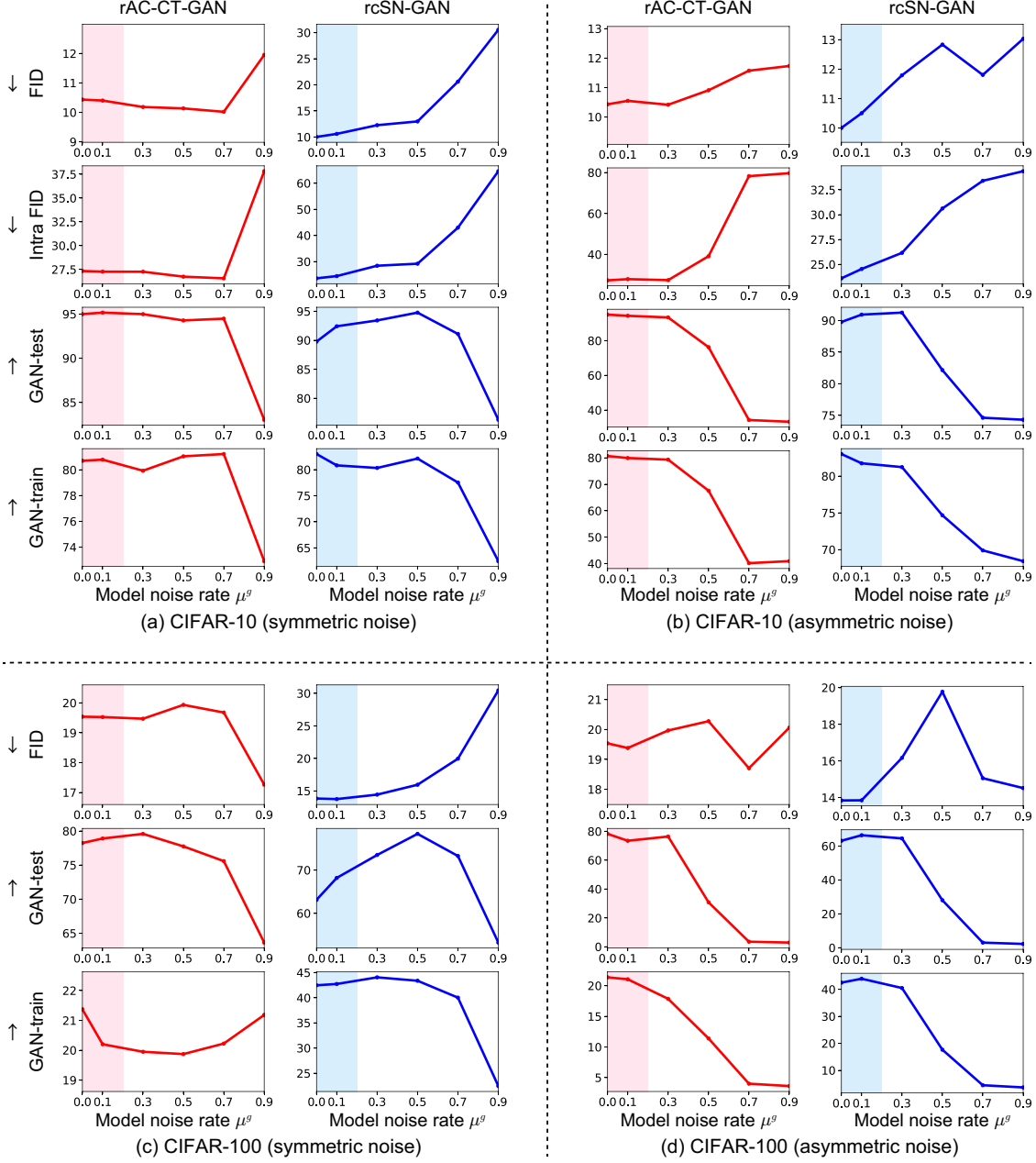
Figure 14. Effect of the gap between real and model $T$. We evaluated rAC-CT-GAN and rcSN-GAN on CIFAR-10 and CIFAR-100 in symmetric and asymmetric noise settings. We fixed a real noise rate as $\mu^r = 0$ and varied a model noise rate $\mu^g$ in $\{0, 0.1, 0.3, 0.5, 0.7, 0.9\}$. The colored area indicates that the gap is within $\pm 0.2$. Note that the scale is adjusted on each graph for easy viewing.
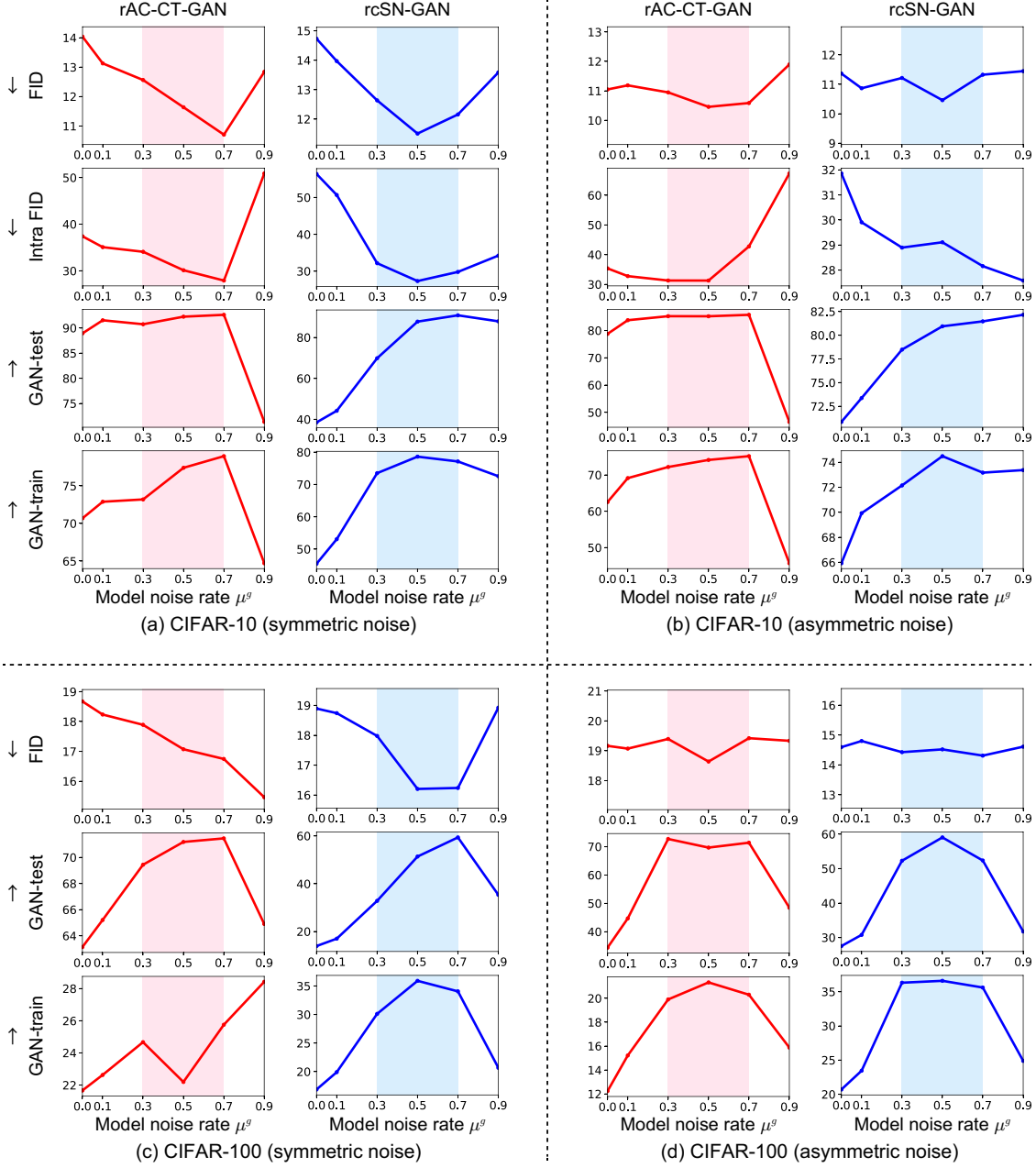
Figure 15. Effect of the gap between real and model $T$. We evaluated rAC-CT-GAN and rcSN-GAN on CIFAR-10 and CIFAR-100 in symmetric and asymmetric noise settings. We fixed a real noise rate as $\mu^r = 0.5$ and varied a model noise rate $\mu^g$ in $\{0, 0.1, 0.3, 0.5, 0.7, 0.9\}$. The colored area indicates that the gap is within $\pm 0.2$. Note that the scale is adjusted on each graph for easy viewing.
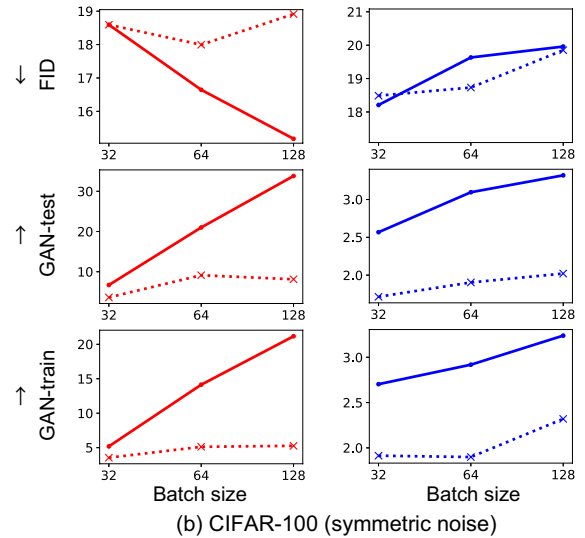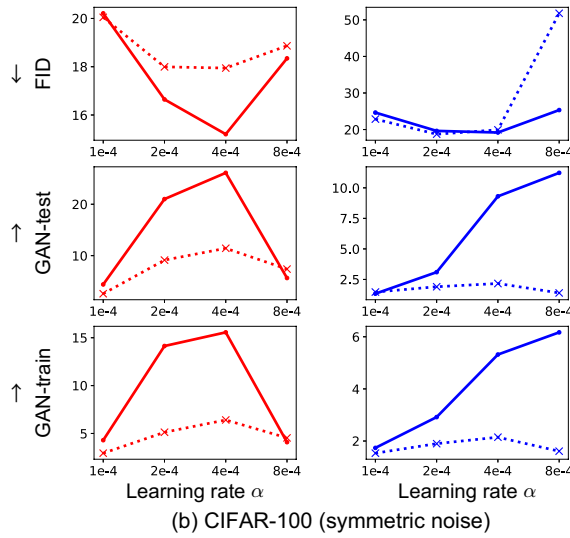
Figure 16. Comparison of the models with learning rates of 0.0001, 0.0002 (default), 0.0004, and 0.0008. We evaluated AC-CT-GAN, rAC-CT-GAN, cSN-GAN, and rcSN-GAN on CIFAR-10 and CIFAR-100 in severely noisy settings (i.e., symmetric noise with a noise rate 0.9). We fixed the batch size as 64 (default). Note that the scale is adjusted on each graph for easy viewing.



Figure 17. Comparison of the models with batch sizes of 32, 64 (default), and 128. We evaluated AC-CT-GAN, rAC-CT-GAN, cSN-GAN, and rcSN-GAN on CIFAR-10 and CIFAR-100 in severely noisy settings (i.e., symmetric noise with a noise rate 0.9). We fixed the learning rate $\alpha$ as 0.0002 (default). Note that the scale is adjusted on each graph for easy viewing.
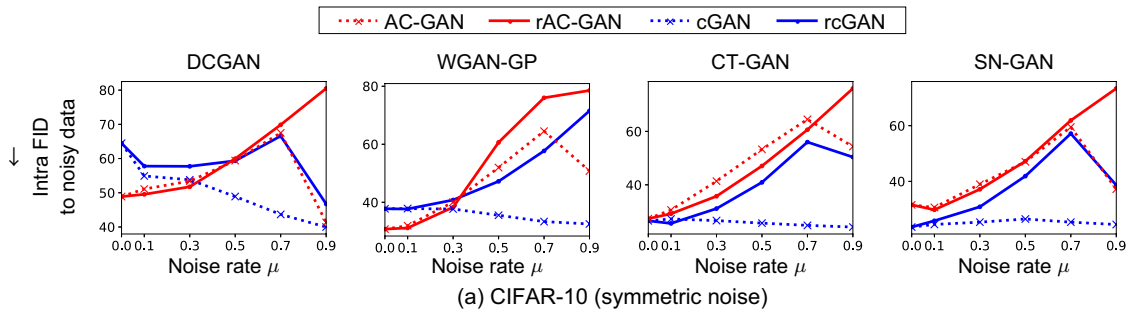
Figure 18. Intra FID between the real noisy data distribution and the generated data distribution. We evaluated AC-GAN, rAC-GAN, cGAN, and rcGAN with combinations of DCGAN, WGAN-GP, CT-GAN, and SN-GAN on CIFAR-10 with symmetric noise. Note that the scale is adjusted on each graph for easy viewing.

# D. Details on Section 7.1

## D.1. Network architectures and training settings

In the experiments on CIFAR-10 and CIFAR-100 (Section 7.1–7.3), we tested four GAN configurations: DC-GAN [54], WGAN-GP [17], CT-GAN [74], and SN-GAN [48]. As discussed in Section 7.1, instead of extensively searching for the best parameters for each label-noise setting, we tested them with the default parameters that are commonly used in clean label settings, and investigated the label-noise effect for them. We explain each one below.

**Notation.** In the description of network architectures, we use the following notation.

- FC: Fully connected layer
- Conv: Convolutional layer
- Deconv: Deconvolutional (i.e., fractionally strided convolutional) layer
- BN: Batch normalization [24]
- ReLU: Rectified unit [50]
- LReLU: Leaky rectified unit [41, 76]
- ResBlock: Residual block [20]
- Concat($y$): Concatenating $y$ ($\in \{1, \ldots, c\}$) after converting it to a one-hot vector ($\in \mathbb{R}^c$) and reshaping it to adjust feature size
- Proj(Embed($y$)): Embedding $y$ such that its dimension becomes the same as of the previous layer $\boldsymbol{h}$ and taking an inner product between embedded $y$ and $\boldsymbol{h}$

In the description of training settings, we use the following notation. Note that we used the Adam optimizer [32] for all GAN training.

- $\alpha$: Learning rate of Adam
- $\beta_1$: The first order momentum parameter of Adam
- $\beta_2$: The second order momentum parameter of Adam
- $n_D$: The number of updates of $D$ per one update of $G$

### D.1.1 DCGAN

DCGAN [54] is a commonly used baseline model. The main principle of DCGAN is to compose the generator and discriminator using only convolutional layers along with batch normalization [24]. It shows promising results in image generation and unsupervised representation learning.

**Network architectures.** We implemented standard CNN network architectures while referring to [48, 52]. We describe their details in Table 8. The conditional generators used in AC-GAN/rAC-GAN and cGAN/rcGAN are the same (Table 8(a)), while the discriminators are different. For AC-GAN/rAC-GAN, we used $D/C$ in which the layers are shared between $D$ and $C$ except for the last layer,

following [52] (Table 8(b)). For cGAN/rcGAN, we used the *concat* discriminator [47] that employs the conditional information by concatenating the conditional vector to the feature vectors (Table 8(c)).

**Training settings.** In DCGAN, a non-saturating loss [16] is used as a GAN objective function. We trained the networks for $100k$ iterations using Adam with $\alpha = 0.0002$, $\beta_1 = 0.5$, $\beta_2 = 0.999$, $n_D = 1$, and batch size of 64. In AC-GAN and rAC-GAN, we set the trade-off parameters $\lambda_{\mathrm{AC}}^r$ and $\lambda_{\mathrm{AC}}^q$ to 1.

| (a) **Conditional generator** $G(\boldsymbol{z}, y)$ |
| --- |
| $\boldsymbol{z} \in \mathbb{R}^{128} \sim \mathcal{N}(0, I)$, Concat($y$) |
| FC $\to 4 \times 4 \times 512$, BN, ReLU |
| $4 \times 4$, stride=2 Deconv 256, BN, ReLU |
| $4 \times 4$, stride=2 Deconv 128, BN, ReLU |
| $4 \times 4$, stride=2 Deconv 64, BN, ReLU |
| $3 \times 3$, stride=1 Conv 3, Tanh |

| (b) **AC-GAN/rAC-GAN discriminator** $D(\boldsymbol{x})/C(\boldsymbol{x})$ |
| --- |
| RGB image $\boldsymbol{x} \in \mathbb{R}^{32 \times 32 \times 3}$ |
| $3 \times 3$, stride=1 Conv 64, BN, LReLU |
| $4 \times 4$, stride=2 Conv 64, BN, LReLU |
| $3 \times 3$, stride=1 Conv 128, BN, LReLU |
| $4 \times 4$, stride=2 Conv 128, BN, LReLU |
| $3 \times 3$, stride=1 Conv 256, BN, LReLU |
| $4 \times 4$, stride=2 Conv 256, BN, LReLU |
| $3 \times 3$, stride=1 Conv 512, BN, LReLU |
| FC $\to 1$ for $D$, FC $\to c$ for $C$ |

| (c) **cGAN/rcGAN discriminator** $D(\boldsymbol{x}, y)$ |
| --- |
| RGB image $\boldsymbol{x} \in \mathbb{R}^{32 \times 32 \times 3}$, Concat($y$) |
| $3 \times 3$, stride=1 Conv 64, BN, LReLU, Concat($y$) |
| $4 \times 4$, stride=2 Conv 64, BN, LReLU, Concat($y$) |
| $3 \times 3$, stride=1 Conv 128, BN, LReLU, Concat($y$) |
| $4 \times 4$, stride=2 Conv 128, BN, LReLU, Concat($y$) |
| $3 \times 3$, stride=1 Conv 256, BN, LReLU, Concat($y$) |
| $4 \times 4$, stride=2 Conv 256, BN, LReLU, Concat($y$) |
| $3 \times 3$, stride=1 Conv 512, BN, LReLU, Concat($y$) |
| FC $\to 1$ |

Table 8. Standard CNN architectures for CIFAR-10 and CIFAR-100. The basic architectures are the same as those in [48]. The slopes of all LReLU are set to 0.1. Following the AC-GAN paper [52], in $D$ we adopted dropout (with a drop rate 0.5) after all convolutional layers.

### D.1.2 WGAN-GP

WGAN-GP [17] is one of the most widely-accepted models in the literature at present. It is an improved variant of WGAN [4] and incorporates a gradient penalty (GP) term as an alternative to weight clipping. By using GP, it imposes a Lipschitz constraint on the discriminator (called the *critic* in that work). This allows for stabilizing the training of a wide

variety of GAN architectures without relying on heavy hyperparameter tuning. We defined the network architectures and training settings based on the source code provided by the authors of WGAN-GP.[7]

**Network architectures.** We used ResNet architectures provided in the WGAN-GP paper [17]. We describe their details in Table 9. As in DCGAN, the conditional generators used in AC-GAN/rAC-GAN and cGAN/rcGAN are the same (Table 9(a)), while the discriminators are different. For AC-GAN/rAC-GAN, we used $D/C$ in which the layers are shared between $D$ and $C$ except for the last layer, following [17] (Table 9(b)). For cGAN/rcGAN, we used the *projection* discriminator [49] that incorporates the conditional information in a projection based manner (Table 9(c)).

**Training settings.** In WGAN-GP, Wasserstein loss and GP are used as a GAN objective function. We set the trade-off parameter between them ($\lambda_{\mathrm{GP}}$) to 10. We trained the networks for $100k$ generator iterations using Adam with $\alpha = 0.0002$ (linearly decayed to 0 over $100k$ iterations), $\beta_1 = 0$, $\beta_2 = 0.9$, $n_D = 5$, and batch size of 64. In AC-GAN and rAC-GAN, we set the trade-off parameters $\lambda_{\mathrm{AC}}^r$ and $\lambda_{\mathrm{AC}}^g$ to 1 and 0.1, respectively.

### D.1.3 CT-GAN

At present, CT-GAN [74] is one of the state-of-the-art models. It is an improved variant of WGAN-GP and adds a consistency term (CT) to impose a Lipschitz constraint for the whole input domain. This contributes a further improvement in stabilizing the training and raises the quality of generated images. We reimplemented the model while referring to the source code provided by the authors of CT-GAN.[8]

**Network architectures.** The network architectures of CT-GAN are the same as those of WGAN-GP except that dropout is used in $D$. We describe its detailed settings in the caption of Table 9.

**Training settings.** The training settings of CT-GAN are the same as those of WGAN-GP except that CT is added in case CT-GAN. We set the trade-off parameter between the Wasserstein loss and CT ($\lambda_{\mathrm{CT}}$) to 2. The other settings (namely, $\alpha$, $\beta_1$, $\beta_2$, $n_D$, $\lambda_{\mathrm{GP}}$, batch size, and number of iterations as well as $\lambda_{\mathrm{AC}}^r$ and $\lambda_{\mathrm{AC}}^g$ in AC-GAN and rAC-GAN) are the same as those of WGAN-GP.

### D.1.4 SN-GAN

SN-GAN [48] is also one of the state-of-the-art models at present. It introduces spectral normalization to impose a

| (a) **Conditional generator** $G(\boldsymbol{z}, y)$ |
| --- |
| $\boldsymbol{z} \in \mathbb{R}^{128} \sim \mathcal{N}(0, I)$ |
| FC $\to 4 \times 4 \times ch$ |
| ResBlock up $ch$ |
| ResBlock up $ch$ |
| ResBlock up $ch$ |
| BN, ReLU |
| $3 \times 3$, stride=1 Conv 3, Tanh |

| (b) **AC-GAN/rAC-GAN discriminator** $D(\boldsymbol{x})/C(\boldsymbol{x})$ |
| --- |
| RGB image $\boldsymbol{x} \in \mathbb{R}^{32 \times 32 \times 3}$ |
| ResBlock down 128 |
| ResBlock down 128 |
| ResBlock 128 |
| ResBlock 128 |
| ReLU |
| Global pooling |
| FC $\to 1$ for $D$, FC $\to c$ for $C$ |

| (c) **cGAN/rcGAN discriminator** $D(\boldsymbol{x}, y)$ |
| --- |
| RGB image $\boldsymbol{x} \in \mathbb{R}^{32 \times 32 \times 3}$ |
| ResBlock down 128 |
| ResBlock down 128 |
| ResBlock 128 |
| ResBlock 128 |
| ReLU |
| Global pooling |
| (FC $\to 1$) + Proj(Embed($y$)) |

Table 9. ResNet architectures for CIFAR-10 and CIFAR-100. The basic network architectures are the same as those in [17, 74, 48]. In $G$'s ResBlock, conditional batch normalization [13, 11] was used to impose a conditional constraint on $G$. Following [17, 74], in WGAN-GP and CT-GAN, we set $ch = 128$ in $G$ and used global mean pooling in $D$. In CT-GAN, we applied dropout (with drop rates of 0.2, 0.5, and 0.5 from the upper block) after the second to fourth ResBlocks in $D$. Following [48], in SN-GAN, we set $ch = 256$ in $G$, used global sum pooling in $D$, and applied spectral normalization to all the layers in $D$.

Lipschitz constraint. This helps stabilizing the discriminator training, and SN-GAN brings a breakthrough in image generation in complex settings (e.g., high-resolution image generation in ImageNet). We reimplemented the model based on the source code provided by the authors of SN-GAN.[9]

**Network architectures.** The network architectures of SN-GAN are the same as those of WGAN-GP except that the feature maps are doubled in $G$, global sum pooling is used instead of global mean pooling in $D$, and spectral normalization is applied to all the layers in $D$. We describe the details in Table 9.

**Training settings.** The training settings of SN-GAN are

---

[7]https://github.com/igul222/improved_wgan_training
[8]https://github.com/biuyq/CT-GAN

[9]https://github.com/pfnet-research/sngan_projection

| GAN | Architecture | $G\ ch$ | $D$ dropout | $D$ pooling | Objective function |
|---|---|---|---|---|---|
| DCGAN | CNN | - | ✓ | - | Non-saturating loss |
| WGAN-GP | ResNet | 128 | ✗ | Mean | Wasserstein loss + gradient penalty (GP) |
| CT-GAN | ResNet | 128 | ✓ | Mean | Wasserstein loss + GP + consistency term (CT) |
| SN-GAN | ResNet | 256 | ✗ | Sum | Hinge loss (with spectral normalization (SN)) |

Table 10. Comparison of network architectures and training settings

also nearly identical to those of WGAN-GP except that a hinge-loss [38, 68] is used instead of Wasserstein loss and GP. The other settings including $\alpha$, $\beta_1$, $\beta_2$, $n_D$, batch size, and number of iterations as well as $\lambda_{\mathrm{AC}}^r$ and $\lambda_{\mathrm{AC}}^g$ in AC-GAN and rAC-GAN, are the same as those of WGAN-GP.

### D.1.5   Summary

We summarize the difference in network architectures and training settings between DCGAN, WGAN-GP, CT-GAN, and SN-GAN in Table 10.

### D.2. Evaluation metrics

As discussed in Section 7.1, we used four metrics for a comprehensive analysis: (1) the Fréchet inception distance (FID), (2) Intra FID, (3) the GAN-test, and (4) the GAN-train. In this appendix, we describe the detailed procedure for calculating the scores.

### D.2.1   FID

The FID [22] measures the 2-Wasserstein distance between $p^r$ and $p^g$, and is defined as

$$
\begin{aligned}
F(p^r, p^g) = &\|\boldsymbol{m}^r - \boldsymbol{m}^g\|_2^2 \\
&+ \mathrm{Tr}(\boldsymbol{C}^r + \boldsymbol{C}^g - 2(\boldsymbol{C}^r\boldsymbol{C}^g)^{1/2}),
\end{aligned} \tag{18}
$$

where $\{\boldsymbol{m}^r, \boldsymbol{C}^r\}$ and $\{\boldsymbol{m}^g, \boldsymbol{C}^g\}$ denote the mean and covariance of the final feature vectors of the Inception model calculated over real and generated samples, respectively. The authors show that the FID has correlation with human judgment and is more resilient to noise or mode-dropping than the Inception score [62] that is also commonly used in this field. We used the FID to assess the quality of an overall generative distribution. In the experiments, we computed the FID between the $50k$ samples generated by $G$ and all the samples in the training set. The implementation was based on the source code provided by the authors of FID.[10]

Generally, in GANs, it is difficult to define the timing when to stop the training, partially because of the lack of an explicit likelihood measure. It is still an open issue, but as an approximate solution, we chose the best model (i.e., simulated early stopping) based on the FID, following [40]. Precisely, we calculated the FID every $5k$ iterations and

chose the best model in terms of the FID. We calculated the other scores (namely, Intra FID, the GAN-test, and the GAN-train) using this model and reported the results in this paper.

### D.2.2   Intra FID

Intra FID [49] is a variant of the FID and calculates the FID for each class. The authors of Intra FID empirically observed that Intra FID had correlation with the diversity and visual quality in conditional image generation tasks. We used Intra FID to check the quality of a conditional generative distribution. In the experiments, we computed Intra FID between the $5k$ samples generated by $G$ and all the samples in the training set belonging to the class of concern. We reported the score averaged over the classes. We used this metric only for CIFAR-10 because in CIFAR-100, the number of clean labeled data for each class is insufficient to calculate Intra FID (which needs to be $\geq$2,048).

### D.2.3   GAN-test

The GAN-test [63] is the accuracy of a classifier trained on real images and is evaluated on the generated images. This metric is developed for conditional generative models and approximates the precision (i.e., image quality) of them. As a classifier, we used PreAct ResNet-18 used in [81], which is an 18-layer network with preactivation residual blocks [21]. The implementation was based on the source code provided by the authors of [81].[11] We used a cross-entropy loss as an objective function and trained 200 epochs with a batch size of 128. We set an initial learning rate to 0.1 and divided it by 10 after 100 and 150 epochs. Weight decay was set to 0.0001. The accuracy scores for the real test sets of CIFAR-10 and CIFAR-100 were 94.8% and 75.9%, respectively (which were the average scores over the last 10 epochs for three classifiers with random initializations). While calculating the GAN-test, we generated $50k$ samples for evaluation. We calculated the accuracy for them using the above three classifiers and reported their average scores.

---

[10]https://github.com/bioinf-jku/TTUR

[11]https://github.com/facebookresearch/mixup-cifar10

### D.2.4 GAN-train

The GAN-train [63] is the accuracy of a classifier trained on generated images and evaluated on real images in a test set. This metric is also developed for conditional generative models and approximates the recall (i.e., diversity) of them. Regarding the classifier, we used the same network architecture and training settings as those used for the GAN-test (described in Appendix D.2.3). While training the classifier, we generated $50k$ samples as training samples. Using this classifier, we calculated the accuracy for the test set. We reported the scores averaged over the last 10 epochs.

## E. Details on Section 7.2

In this appendix, we provide the details of the noise transition probability estimation method used in Section 7.2. As discussed in Section 6.1, we used a *robust two-stage training algorithm* [53], which can estimate $T'$ independently of the main model (namely, an image classification model in [53] and a conditional generative model in our case). In this algorithm, a noisy label classifier $C'(\tilde{y}|\boldsymbol{x})$ is first trained using noisy labeled data and then $T'$ is estimated via the following two steps:

$$\bar{\boldsymbol{x}}^i = \operatorname*{argmax}_{\boldsymbol{x}\in\mathcal{X}'} C'(\tilde{y}=i|\boldsymbol{x}) \tag{19}$$

$$T'_{i,j} = C'(\tilde{y}=j|\bar{\boldsymbol{x}}^i), \tag{20}$$

where $\mathcal{X}'$ is a dataset used for calculating $T'$. In practice, we used the training set as $\mathcal{X}'$. After estimating $T'$, the main model is trained using it.

While implementing the classifier, we used the network architecture and training settings that are similar to those used in calculating the GAN-test (see Appendix D.2.3). However, one possible problem encountered while solving Equations 19 and 20 using DNNs is the memorization effect [80], i.e., $C'(\tilde{y}|\boldsymbol{x})$ can fit to noisy labels and make all probabilities to be zero or one. This causes difficulty in obtaining $T'$ having reasonable values. To alleviate the effect, we added an explicit regularization (i.e., added dropout with a drop rate 0.8 after the first convolutional layer in each residual block) to degrade training performance on noisy labeled data [5], conducted temperature scaling [18] to mitigate the gap between accuracy and confidence, and took a $\alpha$-percentile in place of the $\operatorname{argmax}$ of Equation 19 [45, 53] to eliminate the data strongly fitting the labels. Following [53], we set $\alpha$ empirically. We used $\alpha = 97\%$ for the CIFAR-10 symmetric and asymmetric noise and set $\alpha = 100\%$ (i.e., $\operatorname{argmax}$ is directly used) and $\alpha = 99.7\%$ for the CIFAR-100 symmetric and asymmetric noise, respectively.

## F. Details on Section 7.3

In this appendix, we present the implementation details of the improved technique for severely noise data, which is introduced in Section 6.2 and is evaluated in Section 7.3. Regarding the network architecture, we used shared networks between $D$ (or $D/C$ in AC-GAN/rAC-GAN) and $Q$. Following a sharing scheme between $D$ and $C$ in AC-GAN/rAC-GAN, we shared the layers between $D$ (or $D/C$) and $Q$ except for the last layer. The other parameters that we needed to define were the trade-off parameters between $\mathcal{L}_{\mathrm{GAN}}$ and $\mathcal{L}_{\mathrm{MI}}$, i.e., $\lambda_{\mathrm{MI}}^g$ and $\lambda_{\mathrm{MI}}^q$. We empirically defined the parameters dependent on the GAN configurations (i.e., DCGAN, WGAN-GP, CT-GAN, or SN-GAN) and model (i.e., rAC-GAN or rcGAN) but independent of the datasets (i.e., CIFAR-10 or CIFAR-100). We list them in Table 11.

| Model | GAN | $\lambda_{\mathrm{MI}}^q$ | $\lambda_{\mathrm{MI}}^g$ |
|---|---|---|---|
| Improved rAC-GAN | DCGAN | 0.01 | 0.04 |
| | WGAN-GP | 1 | 0.02 |
| | CT-GAN | 0.01 | 0.04 |
| | SN-GAN | 1 | 0.02 |
| Improved rcGAN | DCGAN | 1 | 0.04 |
| | WGAN-GP | 1 | 0.04 |
| | CT-GAN | 0.01 | 0.04 |
| | SN-GAN | 1 | 0.04 |

Table 11. Hyperparameters for improved rAC-GAN and improved rcGAN

## G. Details on Section 7.4

### G.1. Network architectures and training settings

In the experiments on Clothing1M (Section 7.4), we used two GAN configurations: CT-GAN for AC-GAN/rAC-GAN and SN-GAN for cGAN/rcGAN. We defined the network architectures and training settings while referring to the source code provided by the authors of SN-GAN [48] (which is used for $64 \times 64$ dog and cat image generation).[12] The reason why we refer to this source code is that there is no previous study attempting to learn a generative model using Clothing1M, to the best of our knowledge. We experimentally confirm that its settings are reasonable for Clothing1M with no hyperparameter tuning.

**Network architectures.** We describe the details on the network architectures in Table 12. They are basically similar to those in CIFAR-10 and CIFAR-100 (described in Appendix D.1) except that the input image size is different (that is $32 \times 32 \times 3$ in CIFAR-10 and CIFAR-100, while that is $64 \times 64 \times 3$ in Clothing1M) and feature map size is modified to adjust to the input size difference.

---

[12]https://github.com/pfnet-research/sngan_projection

**Training settings.** In CT-GAN, we set the trade-off parameters to $\lambda_{\mathrm{GP}} = 10$ and $\lambda_{\mathrm{CT}} = 2$, which are the same as those in CIFAR-10 and CIFAR-100. We trained the networks for $150k$ generator iterations using Adam with $\alpha = 0.0002$ (linearly decayed to 0 over the last $50k$ iterations), $\beta_1 = 0$, $\beta_2 = 0.9$, $n_D = 5$, and batch size of 64. We set the trade-off parameters $\lambda_{\mathrm{AC}}^r$ and $\lambda_{\mathrm{AC}}^g$ to 1 and 0.1, respectively. In SN-GAN, we used the same settings except that a GAN objective function is replaced from the Wasserstein loss + GP + CT to the hinge loss.

| (a) **Conditional generator** $G(\boldsymbol{z}, y)$ |
|---|
| $\boldsymbol{z} \in \mathbb{R}^{128} \sim \mathcal{N}(0, I)$ |
| FC $\rightarrow 4 \times 4 \times 1024$ |
| ResBlock up 512 |
| ResBlock up 256 |
| ResBlock up 128 |
| ResBlock up 64 |
| BN, ReLU |
| $3 \times 3$, stride=1 Conv 3, Tanh |

| (b) **AC-GAN/rAC-GAN discriminator** $D(\boldsymbol{x})/C(y)$ |
|---|
| RGB image $\boldsymbol{x} \in \mathbb{R}^{64 \times 64 \times 3}$ |
| ResBlock down 64 |
| ResBlock down 128 |
| ResBlock down 256 |
| ResBlock down 512 |
| ResBlock down 1024 |
| ReLU |
| Global mean pooling |
| FC $\rightarrow 1$ for $D$, FC $\rightarrow$ c for $C$ |

| (c) **cGAN/rcGAN discriminator** $D(\boldsymbol{x}, y)$ |
|---|
| RGB image $\boldsymbol{x} \in \mathbb{R}^{64 \times 64 \times 3}$ |
| ResBlock down 64 |
| ResBlock down 128 |
| ResBlock down 256 |
| ResBlock down 512 |
| ResBlock down 1024 |
| ReLU |
| Global sum pooling |
| (FC $\rightarrow 1$) + Proj(Embed($y$)) |

Table 12. ResNet architectures for Clothing1M. The basic network architectures are defined while referring to [48]. The detailed settings are similar to those described in Table 9. In $G$'s ResBlock conditional batch normalization [13, 11] was used to impose a conditional constraint on $G$. In CT-GAN, we used global mean pooling in $D$ and applied dropout (with drop rates of 0.2, 0.2, 0.5, and 0.5 from the upper block) after the second to fifth ResBlocks in $D$. In SN-GAN, we used global sum pooling in $D$ and applied spectral normalization to all the layers in $D$.

### G.2. Evaluation metrics

As discussed in Section 7.4, we used the FID and the GAN-train as evaluation metrics in these experiments. We

did not use Intra FID because the number of clean labeled data for each class is insufficient to calculate Intra FID. We did not use the GAN-test because Clothing1M is a challenging dataset and we find that a trained classifier tends to be easily deceived by noisy labeled data.

The calculation procedure of the FID is the same as that for CIFAR-10 and CIFAR-100 (described in Appendix D.2.1). We calculated the FID between the $50k$ generated samples and all the samples in the training set (particularly we used $1M$ noisy data). The calculation procedure of the GAN-test is also similar to that for CIFAR-10 and CIFAR-100 (described in Appendix D.2.3). However, we modified the classifier network architecture and training settings so as to obtain the training stability when trained on real clean labeled data. Regarding the network architecture, we used the same network architecture as that for CIFAR-10 and CIFAR-100 except that dropout (with a drop rate 0.5) is used after the first convolutional layer in each residual block. With regards to the training settings, we used a cross-entropy loss as an objective function and trained 200 epochs with a batch size of 128. We set an initial learning rate to 0.01 and divided it by 10 after 100 and 150 epochs. Weight decay was set to 0.01. The accuracy for the real clean labeled test sets was 71.1%.[13] While training the classifier, we generated $50k$ samples as training samples. Using this classifier, we calculated the accuracy for a test set. We reported the scores averaged over the last 10 epochs.

---

[13]This score cannot be directly compared with the scores in the previous studies because we used $64 \times 64$ images but the previous studies used $256 \times 256$ images.