1 We sincerely thank Reviewers and Area Chairs for valuable feedback.

2 **Reviewer 1** • **1. The optimization function for generator:** We really appreciate Reviewer's suggestion. We have
3 checked Reviewer's derivation carefully. We humbly point out that there is a small sign mistake in the derivation:
4 $\lambda_g \Phi^+(G,C^*)$ was used in Eq. 9, but the correct one is $-\lambda_g \Phi^+(G,C^*)$. We have derived generator optimization. Using
5 $\sum_{k=1}^{K} \mathrm{KL}(P_g^{T_k}||P_d^{T_k}) = K \cdot \mathrm{KL}(P_g||P_d)$, and let $\mathcal{V}_\Phi(\mathrm{x}) = \sum_{k=1}^{K} p^{T_k}(\mathrm{x}) \log(\frac{p^{T_k}(\mathrm{x})}{\sum_{k=1}^{K} p^{T_k}(\mathrm{x})})$, we substitute Eq. 11 into Eq. 9:

$$\min_G \mathcal{V}(D^*, C^*, G) = \mathcal{V}(D^*, G) - \lambda_g \left( - K \cdot \mathrm{KL}(P_g||P_d) + \mathbb{E}_{\mathrm{x} \sim P_g^T} \mathcal{V}_\Phi(\mathrm{x}) \right)$$
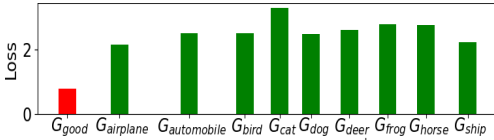
$$= \int_x p_d(\mathrm{x}) \log \frac{p_d(\mathrm{x})}{p_d(\mathrm{x}) + p_q(\mathrm{x})} dx + \int_x p_g(\mathrm{x}) \log \frac{p_g(\mathrm{x})}{p_d(\mathrm{x}) + p_g(\mathrm{x})} dx + K\lambda_g \int_x p_g(\mathrm{x}) \log \frac{p_g(\mathrm{x})}{p_d(\mathrm{x})} dx - \lambda_g \mathbb{E}_{\mathrm{x} \sim P_g^T} \mathcal{V}_\Phi(\mathrm{x})$$

6 Even we assume $K\lambda_g = 1$, it is not easy to simplify $\mathcal{V}(D^*, C^*, G)$. We will further investigate.

7 • 2. Reviewer's comment on analogies to conditional GAN is thought-provoking. Previous work [4] and ours have
8 regarded self-supervised (SS) GAN as unconditional GAN. But Reviewer is correct that SS GAN can also be regarded
9 as conditional GAN with pseudo-labels. Following Reviewer's comment, we will discuss more analogies in the paper.
10 In particular, AC-GAN [30] with pseudo-labels is quite similar to SS-GAN proposed in [4], but deeper analysis is
11 needed to understand the impact of replacing class labels in the original AC-GAN with pseudo-labels.

12 • 3. We will clarify minimax optimization and improve presentation (theorems, typos). We will definitely share code.

13 **Reviewer 2** • **1. Self-supervised (SS) signals:** We thank Reviewer for feedback. Potentially, Contrastive Predictive
14 Coding [*Aaron van den Oord et al. 2018*] can be SS signals for time-series.



Figure A: Our proposed SS loss $-\Phi^+(G,C)$. Red: Good generator; Green: Mode-collapsed generators that correspond to different classes of CIFAR-10, as explained in Sec.4.

• 2. We follow Reviewer's suggestion to use the same toy example to show improvement of our proposed method. The setup is the same as in our paper Sec.4 L167-180, except that now we replace models/cost functions of [4] with our proposed ones, i.e., change from Fig.1a to Fig.1b. Therefore, now the loss is $-\Phi^+(G,C)$, which is shown in Fig.A for a good generator $G_{good}$ and different mode-collapsed generators $G_{collapsed}$ (designs of $G_{good}$ and $G_{collapsed}$ are the same as in our paper Sec.4). Comparing Fig.A and Fig.3a in our paper, the improvement using our proposed model can be observed: $G_{good}$ has
24 the lowest loss under our proposed model.

25 • **3. A separate discriminator to distinguish**
26 **real/fake rotated images:** We train this new model
27 **MS-v2**. Our preliminary result shows that this is very
28 competitive (Table A, row 1). We thank Reviewer for
29 this good suggestion. We will further analyze. We
30 note that this new model further corroborates our idea
31 to leverage discrimination of *rotated* real/fake images
32 to improve *generator* learning, as we have provided
33 theoretical and empirical evidence in our paper.

34 • 4. We show results of CIFAR-100 in Table A, row
35 2. We follow the setup in [R2] exactly, i.e. 10K-5K
36 FID. Note that [R2] has the state-of-the-art results for

| Datasets | SS | MS | MS-v2 |
|---|---|---|---|
| **CIFAR-10** (10K-10K FID) | 12.37 | 11.40 | 11.15 |
| **CIFAR-100** (10K-5K FID) | 49.40 | 21.39 | - |
| **ImageNet $32\times32$** (10K-10K FID) | 26.04 (best) | 13.70 | - |
| **Stacked MNIST** (#mode) | $878.5 \pm 38.9$ | $943.2 \pm 31.4$ | - |
| **Stacked MNIST** (KL) | $0.99 \pm 0.19$ | $0.70 \pm 0.15$ | - |

Table A: Additional results. Baseline: Dist-GAN with ResNet, as in our paper. We also follow exactly the same experiment setup. **SS**: proposed in [4]; **MS**: this work. Note that for **mode coverage** experiment (row 4), our method **MS** achieves the best results for this dataset with tiny $K/4$ architecture, **outperforming state-of-the-art by a significant margin:** [36]: #mode = $859.5 \pm 68.7$, KL = $1.04 \pm 0.29$; [R1]: #mode = $859.5 \pm 36.2$, KL = $1.05 \pm 0.09$. We will update the paper with these results.

37 CIFAR-100. As shown in Table A, row 2, **SS** suffers mode collapse (high FID), probably due to some classes like
38 "sunflowers", "fruit and vegetables", which original SS task does not encourage a generator to learn to produce them
39 (*rotated sunflowers are not rare*). Our **MS** achieves good FID on this dataset and outperforms [R2] (best 10K-5K FID =
40 23.6) on same dataset. Imagenet ($128\times128$ or higher resolution) experiments require very extensive computation. So
41 far, most of Imagenet results are from big companies/institutions, eg., Google [2, 4, 41]. We will look for computing
42 resource for this experiment. We conduct the preliminary experiments on lower resolution Imagenet ($32\times32$) as shown
43 in Table A, row 3. Our **MS** substantially outperforms **SS**. Note that training using our **MS** is also stable: **SS** suffers
44 from mode collapse (FID = 56.20) at the end of training and our **MS** attains the best FID at the end of training.

45 **Reviewer 3** • 1. We will further clarify Fig.3. Note some explanation is in L.167-180. • 2. These methods are compared
46 under **exactly the same settings** of [26, 4]. Some results are missing because previous methods did not report for these
47 experiments. • **3. No. of modes covered**: We thank Reviewer's suggestion. We follow exactly the same experiment
48 setup, tiny architecture $K/4$ and evaluation protocol of [25]. Table A, row 4 shows that our performance for mode
49 covered and KL divergence [25] is superior, significantly outperforms state-of-the-art [36, R1].

50 [R1] Karras et al., Progressive Growing of GANs for Improved Quality, Stability, and Variation, ICLR 2018.
51 [R2] Yamaguchi et al., Distributional Concavity Regularization for GANs, ICLR 2019.