

1 We thank the reviewers for their comments, which we found to be quite helpful. We have strengthened our submission
2 with this feedback, and we believe the below addresses the main concerns identified in our submission.

3 **First** concern is the reviewer believed that we just borrow sparsification algorithms designed for image classification
4 networks, and use them to train generative models. So he thought the baselines are not convincing.

5 We think this is a misunderstanding. For the experiments to compare with our method, we all follow the general settings
6 and schedules when training the generative models, i.e., training in an adversarial fashion, as shown in **Table 1**. Just as
7 comments from other reviewers, we apply the common sparsification algorithms in these baselines is to explain why
8 existing techniques work for models trained on classification tasks but not GANs. So all the baselines are convincing.

9 **Second** concern is the reviewer thought we are missing the clear evidence of some claims about the loss curves.

10 Actually, the detailed evidence of loss curves are already provided in **Appendix, Section A.1**. And the detailed explains
11 are already provided in **Discussion** part in **Section 3** of the manuscript. For example, in which situations, loss curves
12 that look identical to the baseline training will also lead to the bad compression, and the discriminator falls into a
13 low-entropy solution that will cause mode collapse.

14 **Third** concern is the reviewer thought we overload the term “self-supervised”.

15 We thank the reviewers for pointing this out. We think the better term to describe our approach can be “self-tuning”,
16 “self-correcting”, “autoregulative”.

17 **Fourth** concern is the reviewer preferred more mathematical analysis on the performance boost in compression.

18 We provide the analysis from the **Bayes theory** perspective. The three deep neural networks in the **GAN** compression
19 task are the original generative model G_O , the compressed generative model G_C , and the discriminative model D .
20 Given x as the input of the generative networks, we can denote the generative outputs as $G_O(x)$ and $G_C(x)$.

21 We use x_i and x_j to represent two training samples from different categories. Our target is to push closer the generative
22 outputs of the original and compressed generative models with the samples from the same categories, while to push
23 apart the outputs of these two models with the samples from different categories. **KL** divergence is applied to measure
24 the difference between two generative representations. Ideally, the target can be denoted with the following formulas.

$$KL(G_O(x_i), G_C(x_i)) \rightarrow 0, \quad KL(G_O(x_j), G_C(x_j)) \rightarrow \infty \quad (1)$$

25 We define a latent variable S which represents whether the two input samples are from similar ($S = 1$) or different
26 ($S = 0$) categories. For ease of notation, we define the event U to denote the generative representations between the
27 G_O and G_C models are similar, and the event V denotes the D model regards the generative results are similar, i.e.,

$$\begin{aligned} U \Rightarrow G_O(x) \doteq G_C(x), \quad \bar{U} \Rightarrow G_O(x) \neq G_C(x) \\ V \Rightarrow D(G_O(x)) \doteq D(G_C(x)), \quad \bar{V} \Rightarrow D(G_O(x)) \neq D(G_C(x)) \end{aligned} \quad (2)$$

28 According to the total probability formula, for the whole GAN compression process:

$$\begin{aligned} P(S = 1) = P(S = 1 | U, V)P(U, V) + P(S = 1 | \bar{U}, V)P(\bar{U}, V) \\ + P(S = 1 | U, \bar{V})P(U, \bar{V}) + P(S = 1 | \bar{U}, \bar{V})P(\bar{U}, \bar{V}) \end{aligned} \quad (3)$$

29 If the discriminator is initialized by the well-trained model, then the probability of joint distribution for event U and V
30 will be close to $P(U)$, while the probability of joint distribution for \bar{U}, V and U, \bar{V} will be close to 0, simplify (3) as:

$$P(S = 1) = P(S = 1 | U)P(U) + P(S = 1 | \bar{U})P(\bar{U}) \quad (4)$$

31 Because the G_C model is initialized by G_O , so the second item in **formula (4)** has much less influence.

32 If the discriminator is randomly initialized as the original GAN baseline, then U and V can be regarded as the relative
33 independent events. So the four items in **formula (3)** have a certain probability of occurrence.

34 Because the first item in **formula (3)** and **(4)** is our learning target. Our proposed method keeps the same total probability
35 but changes the probability distribution. Because the optimization process during the learning cannot guarantee to find
36 the global optimum. So an easier learning target has a higher expectation to be achieved during the same compression
37 and optimization process. (We will extend to more rigorous prove without the one-page limitation.)

38 **Improvements:** We will add these minor improvements suggested by the reviewers in the final camera-ready version.

39 1. Provide the same level of thorough quantitative and qualitative results comparing to the baseline techniques for at
40 least one more GAN architecture and dataset. (We will add them in **Appendix**.)

41 2. Improve the captions which will be helpful for readers who like to skim figures before deciding to read a paper.

42 3. Add the experiment setting as training a small & dense network from scratch, but with the discriminator initialized
43 as the trained discriminator.