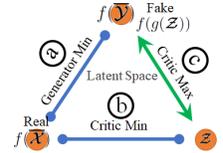


1 **General comments:** We thank all the reviewers for their insightful comments, and their unanimous positive comments
 2 on the superior performance and high quality of our work. Our novelty has also been affirmed by R1, R2 and R4. We
 3 believe that R3 might have misunderstood our work to be an extension of MMD-GANs (e.g., Refs [3,4] suggested by
 4 R3). However, we should clarify that (1) our work differs completely from MMD-GANs, and (2) although Ref [4]
 5 remains unpublished before the submission deadline, our work also achieves superior performances over Refs [3,4]. We
 6 are confident that our work is novel because *our critic operates as semantic embeddings and learns a fruitful latent*
 7 *space*, instead of being a component to build complete metrics as the existing GANs (e.g., MMD-GANs and W-GANs)
 8 do. Theoretically, this avoids the Lipschitz constraint but it requires a reciprocal between the generator and critic, where
 9 the *auto-encoder is a necessity* instead of just a plug-in. Our CF design (even given the unpublished work), is novel in
 10 its • *triangle anchor design with l_1 -norm* (to stabilise convergence), • *meaningful analysis of amplitude and phase* (to
 11 favour other distribution alignment tasks), • *t-net of outputting scales* (to optimise \mathcal{T} distribution types), and • *useful*
 12 *theory*, e.g., Prop. 1 for ill-posed optimisation (a potential problem in Ref [3]) and Lemma 2 for reducing sets. Thus,
 13 our RCF-GAN *seamlessly combines the auto-encoder and GANs by using only two neat modules while achieving s.o.t.a.*
 14 *generation and reconstruction*, whereas MMD-GANs (and many others) typically use three modules, but reconstruct
 15 and interpolate blurred images. Our supplementary material includes the *s.o.t.a.* results under adv-DCGAN and ResNet
 16 structures on more datasets. Below we discuss the reviewers’ comments and will address all of them in the revision.

17 **R1: [Expressivity]:** This is indeed an interesting topic. Although having proved the effectiveness in our RCF-GAN,
 18 the unimodal setting may limit the expressivity on complex datasets and a worthy investigation could be using mixture
 19 models (or learning models) to further improve the performance. **[t-net]:** Yes, exactly. In our code (to be released upon
 20 acceptance) we put both $f(\cdot)$ and t -net in the critic, as they are optimised simultaneously. As t -net is optional (we can
 21 directly use Gaussian samples), we separate it from $f(\cdot)$ but it belongs to the critic; this will be clarified in the revision.

22 **R2: [Q0]:** (1) Fig.2-b shows some results on MNIST by directly training $g(\mathbf{z})$. For other datasets, solely training $g(\mathbf{z})$
 23 via the CF typically performs inferior, which in our preliminary trials on CelebA, obtained a 165 FID score (rough faces).
 24 This also verifies the benefit of latent space comparison via our critic. Tailored t -net could be of help and we will provide
 25 more results in the revision. (2) Yes, other metrics can be seamlessly adopted in our framework. However, we found that
 26 our CF loss works better in practice because our t -net automatically optimises \mathcal{T} distribution types. Our CF loss also
 27 allows for more complex distributions (e.g., mixture models) to further improve model expressibility. **[Q1]:** (1) The
 28 anchor refers to \mathcal{Z} and the critic loss (minimising) is $-(C_{\mathcal{T}}(f(\mathcal{Y}), \mathcal{Z}) - C_{\mathcal{T}}(f(\mathcal{X}), \mathcal{Z}))$. In the dynamic training process,
 29 \mathcal{Z} provides a pivot for the critic; please see the figure. This way, the critic can quickly map real data
 30 \mathcal{X} to the support of \mathcal{Z} via (b) (as Lemma 4 requires). (2) Lemma 4 ensures $\mathbb{E}_{\mathcal{Z}}[\|\mathbf{z} - f(g(\mathbf{z}))\|_2^2] =$
 31 $0 \leftrightarrow \mathbb{E}_{\mathcal{Y}}[\|\bar{\mathbf{y}} - f(g(\bar{\mathbf{y}}))\|_2^2] = 0$, i.e., the equivalence between reconstructing in the latent space (on
 32 generated images) and the pixel domain (on real images). (3) As GANs, we adversarially train
 33 $f(\bar{\mathcal{Y}})$ away from \mathcal{Z} by (c). **[Q2]:** $C_{\mathcal{T}}(\cdot, \cdot)$ is then not a valid metric for the degenerated \mathcal{T} since



34 its support is not \mathbb{R}^m as Lemma 1 requires (may have $C_{\mathcal{T}}(\mathcal{X}, \mathcal{Y}) = 0$ but $\mathcal{X} \neq \mathcal{Y}$). **[Q3]:** We apologise for the typos as
 35 pointed out by the reviewer. For the second question, the reviewer might have missed the need of identical supports on
 36 f and g . If $\mathcal{Z} \in \mathbb{R}^m$, the support of $\mathcal{Y} = (\mathcal{Z}, -1)$ on g is \mathbb{R}^m , which does not equal to the support of (\mathcal{Z}, b) on f (i.e.,
 37 \mathbb{R}^{m+1}), leading to $g(f(\mathcal{Z}, 1)) \neq (\mathcal{Z}, 1)$. **[Q4]:** $(\mathcal{X}, \mathcal{Y})$ are for the random variables in the latent space and $(\bar{\mathcal{X}}, \bar{\mathcal{Y}})$ for
 38 the pixel domain. **[Q8]:** In Fig.2-b, using the amplitude only results in generating “wrong” numbers (“1” for digit 4, “6”
 39 for digit 5), uneven characters, disconnected artefacts, etc, all of which can be relieved by adding the phase in training.

36 **R3: [Novelty]:** Please refer to lines 3-17 of this rebuttal. **[Ablation and anchor design]:** During the rebuttal period
 37 we did the ablation study on CelebA and report the FIDs for the original (15.9), $\lambda = 0$ (59.3) and without anchor (17.5).
 38 Thus, • without reconstruction loss, the generation largely degrades, *validating the correctness and necessity of our*
 39 *reciprocal theory*; • our motivation of anchoring is to find a pivot (i.e., \mathcal{Z}) during dynamic training process, where our
 40 critic can quickly map real data to the support of \mathcal{Z} , as Lemma 4 requires. This is evidenced by the 1.6 FID gain. We
 41 will provide more results in the revision. **[Lipschitz]:** Similar to other GANs (e.g., Fisher GAN and Sphere GAN), the
 42 Lipschitz constraint is not a necessity in our RCF-GAN. Please refer to our proof. **[Consistency]:** The inconsistency on
 43 CelebA with [4] is because [4] rescaled images in ALL datasets to a small size 32×32 . For CIFAR10, our result is
 44 consistent with [4], whereby both our and [4] are inconsistent with [3] (Table 2). Compared to [3,4], our RCF-GAN
 45 also achieved the best performances. **[Notation]:** t_{norm} uses fixed Gaussian and t_{net} learns the t -net for ablation study.

46 **R4: [Model parameters]:** Our critic and generator are evaluated under the almost SAME number of model parameters
 47 as W-GANs, whereas MMD-GANs need an extra decoder net. The only extra cost in our optional t -net is negligible
 48 because it is a 3-layer FC net with the dimension of each layer less than 128. **[Reconstruction and datasets]:**
 49 Fig.4 in the paper shows the image reconstruction and interpolation, validating our superior performances on *clear*
 50 *reconstructions and semantic interpolations*. Our supplementary material further validates our scalable and consistent
 51 superior performances on two extra structures (adv-DCGAN and ResNet) and one extra dataset (LSUN Church). We
 52 will also release our code upon acceptance. **[Phase and amplitude]:** As the amplitude weight measures the model
 53 diversity, the mode collapse can be efficiently relieved in our results. We will elaborate more upon this in the revision.