

1 We thank all reviewers for insightful comments. *All added experiments are tabulated in revised manuscript/appendices.*

2 **[Reviewer 1] • Saito et al.:** We will rename our framework to [TimeGAN](#) to minimize confusion with [\[Saito, ICCV](#)
3 [2017\]](#), which operate *within* the standard GAN framework, proposing a special 2-stage generator (detail in revised
4 appendices). By contrast, we propose a different GAN *framework* altogether, where adversarial learning occurs in the
5 (jointly-optimized) latent space itself. • **Number of variates in experiment data:** While TimeGAN components can
6 indeed be instantiated with various architectures, we focus on the time series setting, using RNNs to illustrate consistent
7 improvement across a variety of data. Media-specific domain applications (e.g. video) are beyond the paper’s scope.
8 However, we agree an even higher-dimensional validation is beneficial. We have conducted [additional experiments](#)
9 on UCI Human Activity Recognition ($dim = 561$). In short, TimeGAN achieves 0.062 (21.2%) & 0.012 (12.7%)
10 discriminative & predictive gains relative to the best benchmark (RCGAN). • **Static vs. temporal features:** We will
11 clarify the following before [lines 106-107](#): "Consider the general data setting where each instance consists of two
12 elements: static features (that do not change over time, e.g. ethnicity, gender), and temporal features (that occur over
13 time, e.g. vital signs, clinical events)." Accommodating static features gives the most general framework, since they
14 often accompany temporal data (e.g. patient data). However, static features are *not* required (we can simply drop the
15 non-recurrent parts of e, r, g, d); the novelties of TimeGAN are in how it handles temporal aspects. • **Further details**
16 **on architecture:** We will publish full source code for TimeGAN with the camera-ready manuscript; this will contain
17 all specifications, training settings, and parameters necessary for reproducibility. Furthermore, in addition to [lines 37-59](#)
18 in the appendices, we will tabulate all technical information with the same granular level of architectural detail (down
19 to dimensions of individual variables) as Appendix B in [\[Lucic, ICML 2019\]](#). Moreover, for additional sensitivities on
20 hyperparameters λ, η , see response [Hyperparameter sensitivities](#) for [Reviewer 3](#). • **Discriminative metric:** To quantify
21 the fidelity of synthetic samples (among other desiderata; see response [Evaluation metrics...](#) for [Reviewer 2](#)), we use the
22 discriminative metric to gauge how indistinguishable samples are from actual data. First, actual sequences are labeled
23 “real”, and sampled sequences “not real”. Then, an off-the-shelf (RNN) classifier is trained to distinguish between the
24 two (a standard supervised task). We are not doing any *pairwise* testing for differences between individual sequences.

25 **[Reviewer 2] • Evaluation metrics, datasets, and benchmarks:** In the familiar application of GANs to images, the
26 vast majority of evaluation relies on inception scores and variants, as well as visual fidelity by inspection; importantly,
27 observe that the former is based on a separately trained model [\[Salimans, NIPS 2016\]](#). This approach is virtually
28 universal [\[Lucic, ICML 2019; Brock, ICLR 2019; Wang, ECCV 2018\]](#); furthermore, using post-hoc classifiers for
29 the evaluation of generative models is well-established [\[Isola, CVPR 2017; Zhang, ECCV 2016\]](#). In the context of
30 time-series GANs, we observe *three* comprehensive desiderata: (1) *fidelity*—samples should be indistinguishable from
31 real data; (2) *diversity*—samples should be distributed to cover the real; and (3) *usefulness*—samples should be just as
32 useful as real data when used for the same predictive purposes (i.e. train-on-synthetic, test-on-real). In our evaluation,
33 the discriminative score, t-SNE/PCA analyses, as well as predictive score respectively give measures of (1), (2), and (3).

34 Our approach to evaluation is *much more comprehensive* than prior works on GANs for time series. First, they do not
35 address (1) and (2) directly. C-RNN-GAN uses a single dataset, relying on hand-crafted measures of audio fidelity.
36 RCGAN uses a post-hoc classifier to evaluate usefulness of samples (i.e. (3)), tested on MNIST and a single real
37 dataset (very low $dim = 4$); other metrics are only applied to synthetic data. By contrast, we focus on [all 3 desiderata](#).
38 Second, we provide experimental results across [6 competing benchmarks](#) over all metrics; RCGAN and C-RNN-GAN
39 provide zero. Third, our [5 datasets](#) are specifically picked to vary with respect to dimensions, correlations, periodicity,
40 discreteness, etc. (see [lines 239-254](#)), including a massive ($n = 150k, dim = 54$) real-world medical dataset (Events).
41 For these reasons, we submit that our approach to evaluation is *more comprehensive*—especially w.r.t. RCGAN as the
42 reviewer mentions. Furthermore, see response [Number of variates in experiment data](#) for [Reviewer 1](#) for additional
43 results on even higher-dimensional ($dim = 561$) data. • **Static vs. temporal features:** Kindly refer to response [Static](#)
44 [vs. temporal features](#) for [Reviewer 1](#). • **Discrete data:** We already use discrete data. Our largest real-world dataset
45 ($dim = 54$) is discrete, with $\sim 150k$ sequences (see [lines 252-254](#), and [Table 2](#) in appendices). TimeGAN significantly
46 outperforms all benchmarks on both discriminative/predictive scores (as with all datasets. See [Table 2](#) in manuscript).

47 **[Reviewer 3] • Difficulty of training:** Although GANs in general are not the easiest to train, we did not discover any
48 additional complications in our experiments. The embedding task serves to regularize adversarial learning—which
49 now occurs in a lower-dimensional latent space; similarly, the supervised loss has a constraining effect on the stepwise
50 dynamics of the generator. For both reasons, we do not expect TimeGAN to be *more* challenging to train; standard
51 techniques for improving GAN training still apply. Here, we use [covariance feature matching](#) across all models to
52 improve the diversity of generation. Kindly refer to response [Number of variates in experiment data](#) for [Reviewer 1](#)
53 for equally favorable results on even higher-dimensional data. See also the following response. • **Hyperparameter**
54 **sensitivities:** We find empirically that TGAN is not very sensitive to λ and η . While we set $\lambda = 1, \eta = 10$ for all
55 experiments, we agree that showing [additional sensitivities](#) may be beneficial. We have conducted experiments across a
56 range of hyperparameters, tabulated in the revised appendices. For example (for Stocks), across $\lambda = \{1, 5, 10, 20\}$ and
57 $\eta = \{0.1, 0.5, 1, 2, 5\}$, the min and max discriminative scores are 0.097 and 0.108, with variance 0.004—showing that
58 performance is by no means brittle—thereby providing further reassurance that TimeGAN is not more difficult to train.