

1 We thank all reviewers for their input on our paper. All reviewers recognized the novelty and appeal of the repulsive
 2 term in our proposed loss, which is indeed the main technical contribution of this work. Since many previous works use
 3 spectrogram-based losses without such a term, and since adding the term greatly boosts performance, we also expect
 4 this contribution to have the greatest impact on future research.
 5 The main improvement asked for by the reviewers is the addition of a set of ablation experiments to more clearly tease
 6 apart the contributions of the different loss terms and modeling choices. We’re currently hard at work running and MOS
 7 testing all mentioned combinations of losses and models, and we report currently obtained results below. All results for
 8 all experiments proposed by the reviewers will be added to the final version of our paper.

9 **New ablations (R3, R4)** Partial results for requested new ablations. Missing MOS scores were not yet finished and
 10 will be added later.

model→	r+m+G	r+m	r+256+G	r+512+G	r+256	r+512	
cFSDS→	0.033	0.033	still running	still running	0.035	0.034	
MOS→	4.25 ± 0.06	4.10 ± 0.07			3.44 ± 0.07	2.89 ± 0.09	
model→	m+G	m	256+G	512+G	256	512	G
cFSDS→	0.039	0.039	0.200	0.047	0.040	0.038	0.075
MOS→		3.00 ± 0.07					4.16 ± 0.06

Table 1: Showing results for all combinations of (1) repulsive term (**r**) yes/no, (2) multi-scale (**m**) or single window size (256/512) or no spectrogram loss, (3) GAN loss (**G**) yes/no. Note that these ablations sampled the cFSDS validation set uniformly, where we used length-weighted sampling for the submitted paper and the table below.

model→	GED + full GAN-TTS	GED + uncond. GAN	GAN-TTS only
cFSDS→	still running	0.040	0.077
MOS→		4.25 ± 0.06	4.16 ± 0.06

Table 2: Results for combining our proposed GED loss with full GAN-TTS, including the conditional discriminators (results pending), and comparing against GED + unconditional GAN, and GAN-TTS, as requested by **R4**.

11 **Small flow-based models do exist (R3, R4)** Our related work section now acknowledges that WaveFlow has made
 12 great progress in reducing the size of flow-based models while maintaining generation quality. A remaining advantage
 13 of our model is that it is still more parallel / has fewer sequential steps, as our models are fully convolutional while
 14 WaveFlow is a hybrid between autoregressive and parallel flow-based models.

15 **Limited performance gains against baseline (R4)** Our model with only the GED loss performs about the same as
 16 the SOTA GAN-TTS model from Binkowski et al. (MOS scores are not statistically different). However, our method
 17 is easier to train and converges faster. GED (which includes the repulsive term) also dramatically improves upon the
 18 baseline without repulsive term: Since using spectrogram-based losses without repulsive term is standard practice, we
 19 feel that comparison against this baseline is most informative in forecasting how useful the proposed techniques will be
 20 for the wider community.

21 **Repulsive term not highlighted or studied enough (R4)** Reviewer 4 notes that the repulsive term in our proposed
 22 loss is the main technical novelty in the paper, but that it should be highlighted and studied more. We now put even
 23 more emphasis on this contribution, and we add new ablation studies to better understand its impact (see above).

24 **GED+iSTFT benefits from directly generating spectrograms (R4)** Our iSTFT generator indeed has an inductive
 25 bias that might make it easier to do well on the cFSDS metric, however it also does very well on MOS. To avoid any
 26 possible confusion: the model still generates raw waveform audio, not just spectrograms.

27 **Relation to Improved Techniques for Training GANs (R1)** We now discuss the relationship between this paper and
 28 our proposed technique in the related work section. In short: the attractive term is indeed similar to the feature matching
 29 term in this paper if one would replace the discriminator activations of the feature matching loss with spectrogram
 30 representations. Our repulsive term directly maximizes the distance in feature space between samples, whereas mini
 31 batch discrimination injects sample dissimilarity within a mini batch as side information into the discriminator which is
 32 trained with the usual discriminator loss of a GAN.

33 **Add definition of $p(x), q(y)$ after equation 1. (R3)** Done.

34 **New citation A new framework for CNN-based speech enhancement in the time domain. (R4)** Added to related
 35 work section.

36 **More information on how to produce the linguistic features (R3)** Unfortunately the linguistic features that are used
 37 as conditioning input for our speech generation model, and which are generated from the source text by a separate
 38 model, are indeed not reproducible using publically available code. We now more fully describe these features in the
 39 appendix and, for comparison, we reference many papers that have used these features before.