

1 We thank our reviewers for their time and valuable comments.

2 **Motivation** We have observed in the literature and also from personal communication at recent conferences incl. ICML
3 and ICLR that almost all text GAN practitioners do not believe it is possible to train a GAN using REINFORCE
4 on language with such high dimensional action spaces (e.g. vocabulary sizes of 10,000 or 20,000 as we have done).
5 Instead the area publishes and promotes complex training techniques to overcome instability, using pre-training of the
6 discriminator and generator; or periodic teacher forcing with a tuned schedule of regularity. Furthermore in some cases
7 the GAN models are heavily pre-trained and only fine-tuned as a GAN with a miniscule learning rate (e.g. "Adversarial
8 Feature Matching for Text Generation" Zhang et al. 2017).

9 We feel this paper will have a significant impact, by showing that stable training can be obtained with REINFORCE.
10 We think this will re-focus the community from overcoming stability to benchmarking with richer data (EMNLP and
11 WikiText are probably too small) and scaling with larger models — hopefully to a state where one can observe a
12 significant difference in sample quality.

13 **Overstatement** The reviewers note that some of the wording in the paper suggests that we have solved GANs for text,
14 which we agree is not the case. We do show ScratchGAN is producing samples of similar quality to language models
15 for these datasets, however it is clearly a much worse generative model than the pure MLE variant and so far, the more
16 compute-intensive GAN training is not providing us with a much better model. We will tone down any language which
17 suggests ScratchGAN outperforms MLE. But we stand by the observation that sample quality and diversity appears to
18 be close to the MLE model, and some metrics confirm this (BLEU, FED).

19 **Sample quality** We certainly agree that neither the MLE or ScratchGAN are producing groundbreaking samples -
20 and we mostly attribute this to choice of dataset. We ran this on two LM datasets that had been benchmarked by
21 prior GAN work (EMNLP and WikiText). This had the benefit of comparison to prior work for objective measures
22 BLEU/self-BLEU (Fig 2a). However it has the downside that the samples are quite bad and make for a difficult
23 qualitative comparison. We think future work should focus on scaling to larger datasets, generating larger bodies of text
24 and using aggregate human evaluation to provide a more objective sense of sample quality.

25 **Reviewer 1** We agree with your point that we are dismissing non-autoregressive language models. We will add a couple
26 of sentences to highlight progress in feed-forward approaches.

27 We have addressed these typos, thank you for noting them! Also, we have increased the tables 7 and 9 from 5 samples
28 per (model, dataset) to 15.

29 **Reviewer 2** Since there is a symbiosis in ScratchGAN between the discriminator and generator — both are recurrent
30 models over text, we do believe a promising direction would be to scale the dataset and model (e.g. to a transformerxl)
31 to obtain better quality samples.

32 We agree MLEs are still superior to GANs, we have not solved GANs for text but we think this work is an important
33 data point along the path to doing so. Please see our *Overstatement* section — in short we will tone down the claims of
34 the paper. Re. code release, we are in the process of trying to release a simple colab script for training, such that people
35 can see all of the components working.

36 **Reviewer 3** We agree the sample quality is not very good, we partly address this in the *Sample quality* section above.
37 It was not clear to us that there was a qualitative difference between MLE and ScratchGAN, some MLE samples are
38 quite degenerative, e.g. “after the sets of UNK wear UNK and UNK ’ UNK ’ UNK to tell him , UNK UNK they play
39 UNK UNK with UNK around a UNK .” However future work should benchmark these approaches at scale, with a larger
40 model and use a cumulative human evaluation to assess qualitative appearance; alongside the automatic scores.

41 For objective measures, we do compare to existing GAN approaches (Fig 2a). However the real objective is to have
42 GANs considerably outperform MLE, since it is agreed this is still the best approach for text generation.

43 Re. overclaiming, we agree ScratchGAN does not outperform MLE and have toned down any language that appears to
44 make this claim (see *Overclaim* section above). We genuinely do not want to claim that ScratchGAN solves GANs for
45 text, just that it is possible to train a GAN to a decent level of quality (judged by objective measures) without a complex
46 training procedure of pre-training, teacher forcing, Gumbel Softmax with a scheduled temperature increase etc.

47 Good point re. pre-trained word embeddings. We decided to keep them because there was no change in performance
48 within the ablation study — for this comparative run the model was truly from-scratch ;-).

49 Aside from language applications, we think the result of this study — that REINFORCE can be stably trained in
50 this challenging setting — will be of interest also to reinforcement learning practitioners that are interested in high-
51 dimensional action spaces; e.g. for medical treatment prediction in electronic health record time-series. Thus we ask
52 you to consider this core research contribution, when reconsidering your score.