

1 Thanks all the reviewers for the comments and suggestions!

2 To Reviewer #1

- 3 • **About 2-head attention** Compared with the attention alignments in machine translation, the attention in TTS
4 is monotonous and much simpler than machine translation. Therefore, we choose fewer attention heads. Our
5 preliminary experiments also show that more heads do not bring much difference to Transformer TTS. We will
6 report the numbers in the new version.
- 7 • **Generated v.s. Groundtruth mel-spectrogram** Knowledge distillation is widely used in non-autoregressive
8 machine translation [1] and speech synthesis [2] for transfer knowledge from the autoregressive teacher model
9 to the non-autoregressive student model. The intuitive explanation is that the teacher model can generate data
10 with smoother distribution and less noise, which is easy to be fitted by the student model [1][2]. The ablation
11 study in Table 4 demonstrates the effectiveness of knowledge distillation.
- 12 • **Unmatched hyperparams** Our FastSpeech model differs from the original Transformer TTS model in that we
13 use Conv1D instead of the original dense network after multi-head attention. So in our experiment, FastSpeech
14 model is initialized from the teacher model with the same configuration, but not the original Transformer TTS
15 model as shown in Appendix A. We will explicitly point out this in the new version of the paper.
- 16 • **Missing hyperparams in Appendix** Thanks for your reminder. We will add the missing hyperparams of
17 pre-net and post-net in the new version of paper. For Transformer TTS, the hidden dimension and the number
18 of layers of CNN are 512, 5 for post-net, and are 512, 3 for pre-net. For FastSpeech, the hidden dimension of
19 the decoder's final linear layer is 80.
- 20 • **Batch size for inference** The batch size is 1 for inference evaluation, in order to simulate the scenario of
21 online production. Many previous works (e.g., [1][2][3]) use this batch size to evaluate the inference latency.
- 22 • **About the inconsistent latency numbers** There is a typo in Table 2. The unit of latency should be "second".
23 We will fix it in the new version of the paper.
- 24 • **Robustness test for Tacotron 2** We evaluate the robustness of Tacotron 2 on the 50 hard sentences. The
25 repeating, skipping and error sentences are 4, 11 and 12 respectively, and the error rate is 24%. We will add
26 the results in the new version of the paper.
- 27 • **The reference for CMOS evaluation** We will add the reference for CMOS in the new version of the paper.

28 [1] Gu, Jiatao, et al. "Non-autoregressive neural machine translation." ICLR 2018.

29 [2] Oord, Aaron van den, et al. "Parallel wavenet: Fast high-fidelity speech synthesis." ICML 2018.

30 [3] Prenger, Ryan, et al. "Waveglow: A flow-based generative network for speech synthesis." ICASSP 2019.

31 To Reviewer #2

- 32 • **About pronunciation dictionary** Our grapheme-to-phoneme tool works like this: it first looks up the pro-
33 nunciation dictionary, and if the dictionary does not contain the word, it predicts the phoneme using a
34 grapheme-to-phoneme model.
- 35 • **Enforce monotonicity** Thanks for your suggestion. We observe from the experiments that there are roughly
36 two kinds of attention in Transformer TTS: diagonal and non-diagonal. If we enforce all the attentions to be
37 diagonal, the performance of Transformer TTS will drop, because the non-diagonal attentions are also very
38 helpful to the model. So we select the diagonal attentions rather than enforce all attentions to be monotonic.
- 39 • **Add the citation** We will add this citation in the new version of the paper.

40 To Reviewer #3

- 41 • **About the hard alignments** We've tried to add Gaussian smoothing to soften the inputs of the FFT block in
42 the mel side. However, it doesn't improve the performance. We think the Conv1D and self-attention can mix
43 the information from neighbors and play a role in softening the hard copied hidden states.
- 44 • **About diagonal alignments** We visualize the attention of each head in Transformer TTS and find that nearly
45 half of the attentions are always diagonal for each data pair. We also find that if an attention has high focus
46 rate, it is always diagonal. We find that some errors are closely related to the poor attention: jumping attention
47 causes skipping and non-monotonic attention causes repeating. We think that diagonal attentions can ensure
48 correct alignments in transformer TTS. We could not find any cases that have no diagonal attention at all.
- 49 • **About the title** Thanks for your advice! We will change the title accordingly.