

1 Thanks all the reviewers for the detailed and thoughtful comments.

2 **(R1, R2) Resemblance to HMM-based speech synthesizers and recognizers** We found resemblance to the previous
3 HMM-based works [1, 2, 3], all of which proposed methods to estimate alignments from unsegmented data. Thanks
4 for commenting the missing references, we will mention about the previous works in the new version of the paper.
5 However, we claim that the contribution of our proposed method is not only about the efficient alignment search, but
6 also about enabling parallel sampling over the complex data distribution. We believe incorporating bijective flows for
7 parallel sampling does not naturally come out from HMM-based works, as HMMs are inherently sequential models.
8 Likewise, estimating alignments in the latent space through dynamic programming is distinctive from the previous
9 parallel TTS works such as FastSpeech. Finally, compared to HMM-VAE [3], Glow-TTS differs in that 1) it focuses
10 on generating samples in parallel not estimating sequences of latent discrete variables, 2) does not need an additional
11 encoder to approximate the latent posterior nor 3) independence assumption of outputs across time.

12 **(R1, R3) Deterministic duration prediction** Our main goal was to build a parallel TTS model which jointly learns to
13 model data distribution and align by itself in an efficient way, and the key components of Glow-TTS were the MAS and
14 the flow-based decoder. We've not thoroughly explored to improve the duration predictor and simply follow the same
15 architecture of FastSpeech. But as the feedback of Reviewer 1, exploring a joint model of durations would be promising
16 in the future direction of our work and we expect it could be effective to model diverse speech.

17 **(R1, R3) Grouped 1x1 Convolutions** We design the grouped 1x1 convolutions to be able to mix channels. For each
18 group, as Reviewer 3 commented, the same number of channels are extracted from one half of the feature map separated
19 by coupling layers and the other half, respectively. For example, if a coupling layer divides a 8-channel feature map (a,
20 b, c, d, e, f, g, h) into two halves (a, b, c, d) and (e, f, g, h), we implement to group them into (a, b, e, f) and (c, d, g, h)
21 when the number of groups is 2. We will add the implementation details and update Figure 8c more clearer.

22 **(R1) Importance of parallel TTS** With parallel vocoders such as WaveGlow, the primary computational bottleneck
23 in mel spectrogram-based end-to-end TTS systems became TTS models. For example, to generate a speech of 5.8
24 seconds on a GPU, WaveGlow spends 120 ms while Tacotron 2 and Glow-TTS require 640 ms and 40 ms, respectively.
25 Therefore, adopting parallel TTS models significantly improves the sampling speed of end-to-end systems.

26 **(R3, R4) Controllability and diversity of sampling** We agree Glow-TTS shows restricted controllability and diversity
27 of some properties of speech such as pitch or intonation as Reviewer 3 pointed out. But, as Glow-TTS estimates the
28 distribution of speech given text tokens, it has the potential to generate diverse speech and control the characteristics of
29 speech, which is not possible with the deterministic models like FastSpeech. The stochasticity comes from the prior
30 distribution of the latent variable z , and we believe that synthesizing with careful sampling from the latent space enables
31 us to control the characteristics of generated samples. In Section 5.3, we showed that varying temperature can change
32 the pitch of generated samples and different z s correspond to speech samples with different intonations. To clarify
33 our analysis, we will add more samples for a variety of text samples in our demo page. In this work, we provided the
34 limited analysis on the controllability of Glow-TTS. Further analysis on how the latent space of z is related to the
35 characteristics of generated samples would be helpful for emotional control, and is left as future work.

36 **(R3) The iterative objective** The reason why we demonstrate the iterative objective provides a good lower bound of
37 the global solution is that the iterative procedure is actually one example of widely used Viterbi training [2], which
38 maximizes log likelihood of the most likely hidden alignment. We also empirically showed that our system works well
39 with the objective. We will add a reference about Viterbi training.

40 **(R3) CMOS results** We conducted CMOS evaluation between Tacotron 2 and Glow-TTS with the sampling temperature
41 0.333. Through 500 ratings on 50 items, Glow-TTS wins Tacotron 2 by a gap of 0.934, which shows preference towards
42 our model over Tacotron 2. We will add this result in the new version of the paper.

43 **(R2, R3) Multi-speaker MOS evaluations** We will add comparison to Tacotron 2 with speaker embeddings as a
44 baseline in the new version of the paper. As for methodology of MOS evaluations, we randomly sampled 50 sentence
45 and speaker id pairs from the test set, which resulted in 40 unique speakers (one sentence for 31 speakers, two sentences
46 for 8 speakers, and three sentences for one speaker). When compared to Tacotron 2 with speaker embeddings, we will
47 sample one sentence for each speaker for the evaluation. We will also include the test phrases as well as speaker ids
48 used for all MOS evaluation in the supplemental material.

49 [1] Tokuda, Keiichi, et al. "Speech synthesis based on hidden Markov models." Proceedings of the IEEE 101.5 (2013): 1234-1252.

50 [2] Rabiner, Lawrence R. "A tutorial on hidden Markov models and selected applications in speech recognition." Proceedings of the
51 IEEE 77.2 (1989): 257-286.

52 [3] Ebberts, Janek, et al. "Hidden Markov Model Variational Autoencoder for Acoustic Unit Discovery." INTERSPEECH. 2017.