Human-in-the-loop Speaker Adaptation for DNN-based Multi-speaker TTS

RESEARCH & DEVELOPMENT PROGRAM

Conventional Method

Transfer-learning-based method [1] extracts sp embedding from reference audio with speaker pretrained in the speaker verification task.

Issue: The conventional method cannot be used reference speech of the target speaker.



Proposed Method

 Target speaker embedding is explored with Sequences Search (SLS) [2] algorithm.

Reference speech is not fed to speaker encod

- Two measures are taken to search for a more na
 - 1. Initialization with mean male/female speaker
- 2. Setting the search space to a quantile of the train



[1] Jia et al., Proc. NeurIPS, 2018. [2] Koyama et al., ACM Trans. on Graphics, 2017. [3] Maekawa, Proc. ICLR, 2021. [5] Takamichi et al., AST, 2020. [6] Kong et al., Proc. NeurIPS, 2018.

Kenta Udagawa, Yuki Saito, Hiroshi Saruwatari

The University of Tokyo, Japan.

Introduction

Speaker adaptation of speech synthesis is a technique for synthesizing speech of unseen speaker with small amounts of data. The conventional speaker adaptation method cannot be used without the reference speech of the unseen speaker. We propose a speaker adaptation method which incorporates humans' perceptual feedback and does not feed reference speech to speaker encoder. Experimental results indicate that the proposed method can synthesize speech with the same or better quality than the conventional method, depending on the user and test speaker.

	Experimental Evaluations			
beaker Sencoder	Experimental Condition			Objective/Sul
	Corpus (speaker encoder	The Corpus of Sp) (947 men and 470 v	oontaneous Japanese (CSJ) [3] women, 660 hours)	(Multiple
without the	TTS model Corpus	FastSpeech 2 [4] The parallel data	of the Japanese Versatile Speed	<pre>ch Comparison methods: SLS-{best, mean, worst}:</pre>
	(TTS model)	(JVS) corpus [5] (49 men and 51 wo speaker)	men, 22 hours 100 sentences per	 spectrogram MAE was {n TL: Transfer-learning-bas Mean-Speaker: Endpoint
edding of et Speaker	Data train splitting test	90 speakers (44 men, 46 women) 4 speakers (2 men, 2 women)		Mel spectrogram MAE res
	val Vocoder	6 speakers (3 mer Pretrained univer	n, 3 women) rsal model of HiFi-GAN [6]	 SLS-worst were inferior to In several cases, SLS-be
	Objectiv	e Evaluatio	on (During Search	jvs078 (male) jvs005 (male)
uential Line	We calculat method is n	ed mel spectrogram nanipulated by 8 pa	m MAE at each step when ou articipants up to 30 steps.	
der. atural voice. ning data	 Methods: Proposed: Proposed method TL: Transfer-learning-based method [1] Mean-Speaker: Endpoint of initial line segment Results: 			 Mean-Speaker TL Mean-Speaker TL The proposed method color many cases.
Speaker nbedding space	 Mel spe manipu The pro compa 	ectrogram MAE ter Ilation. oposed method courses rable to TL depend jvs078 (male)	nded not to improve by huma uld synthesize speech ling on the participant.	$n \qquad jvs078 \text{ (male)} \qquad jvs005 \text{ (male)} \\ \hline \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ $
g en Waveform	1.1 1.0 1.0 0.9 0.9 0.9 0.9 0 0.9 0 0.9 0 0.9 0 0.9 0 0.9 0 0.9 0 0.9 0 0.9 0 0.9 0 0.9 0 0 0 0 0 0 0 0 0 0 0 0 0	5 10 15 20 25 30 jvs005 (male) 5 10 15 20 25 30 SLS	1.2 1.1 1.0 0.9 0.8 0 5 10 15 20 25 30 jvs010 (female) 1.4 1.3 1.2 1.1 1.0 0 5 10 15 20 25 30 jvs010 (female) 1.4 1.3 1.2 1.1 1.0 5 10 15 20 25 30 Jvs010 (female) TL — Mean-Speaker	Ground-Truth Mean-Speaker Similarity MOS results: • In some cases, the proposing of speaker adaptate jvs078 (male) jvs005 (male) • Jvs005 (male) Jvs005 (male)
				Mean-Speaker TL

- Speech samples page:

AVATAR SYMBIOTIC SOCIETY

http://sython.org/demo/udagawa22interspeech/demo_slstts.html



bjective Evaluation e Utterances)

The speaker embedding whose mel ninimum, mean, maximum} sed method [1] of initial line segment

sults:



uld synthesize speech as natural as TL in



bsed method was inferior to Mean-Speaker. sed method could achieve the same tion as TL.

