

# Sample Efficient Adaptive Text-to-Speech

Yutian Chen, Yannis Assael, Brendan Shillingford, David Budden, Scott Reed, Heiga Zen, Quan Wang, Luis C. Cobo, Andrew Trask, Ben Laurie, Caglar Gulcehre, Aäron van den Oord, Oriol Vinyals, Nando de Freitas



DeepMind & Google

## Highlights

- Solving many tasks with few data, as opposed to solving few tasks with many data.
- Few-shot meta-learning enables us to **imitate a new voice-style** with 5 minutes of data as opposed to the previous requirement of tens of hours, and without compromising on quality.
- Our models contain both task-dependent and task-independent cores. This division facilitates training and enables fast adaptation to novel voices with low sample complexity.

Naturalness (Mean Opinion Score)

Subjects were asked to rate the naturalness of generated utterances on a five-point Likert Scale (1: Bad, 2: Poor, 3: Fair, 4: Good, 5: Excellent).

Dataset	LibriSpeech		VCTK		
Real utterance	$4.38\pm0.04$		$4.45\pm0.04$		
van den Oord et al. (2016)	$4.21\pm0.081$				
Nachmani et al. (2018)	$2.53 \pm 1.11$		$3.66\pm0.84$		
Arik et al. (2018)					
adapt embedding	-		$2.67\pm0.10$		
adapt whole-model	-		$3.16\pm0.09$		
encoding + fine-tuning	-		$2.99\pm0.12$		
Jia et al. (2018)					
trained on LibriSpeech	$4.12 \pm 0.05$		$4.01 \pm 0.06$		
Adaptation data size	10s	<5m	10s	<10m	
SEA-ALL (ours)	$3.94\pm0.08$	$4.13 \pm 0.06$	$3.92 \pm 0.07$	$3.92\pm0.07$	
SEA-EMB (ours)	$3.86\pm0.07$	$3.95\pm0.07$	$3.81\pm0.07$	$3.82\pm0.07$	
SEA-ENC (ours)	$3.61\pm0.06$	$3.56\pm0.06$	$3.65\pm0.06$	$3.58 \pm 0.06$	

- We achieve state-of-the-art in both sample naturalness and **voice similarity** with merely a few minutes of audio data.



#### **Training and Adaptation**

The adaptation consists of three stages:



### Voice Similarity (Mean Opinion Score and d-vectors)

(MOS) Subjects were asked to rate the similarity between real and synthesized samples of the same speaker. (1: Not at all similar, 2: Slightly similar, 3: Moderately similar, 4: Very similar, 5: Extremely similar)

Dataset	LibriS	Speech	VCTK	
Real utterance	$4.30\pm0.08$		$4.59\pm0.06$	
Jia et al. (2018)				
trained on LibriSpeech	$3.03\pm0.09$		$2.77\pm0.08$	
Adaptation data size	10s	<5 <i>m</i>	10s	<10m
SEA-ALL (ours)	$3.41 \pm 0.10$	$3.75 \pm 0.09$	$3.51 \pm 0.10$	$3.97 \pm 0.09$
SEA-EMB (ours)	$3.42\pm0.10$	$3.56\pm0.10$	$3.07\pm0.10$	$3.18\pm0.10$
SEA-ENC (ours)	$2.47\pm0.09$	$2.59\pm0.09$	$2.07\pm0.08$	$2.19\pm0.09$

#### **Detecting Synthesized Samples**

- Preliminary results show that a linear SVM trained on the d-vectors of real and synthesized multi-speaker samples from our model achieves 85% accuracy.

## a) Training: Train a multi-speaker model with a shared WaveNet core and independent learned embeddings for each speaker (task)



b) Adaptation: Adapt to a new speaker with few-shot data

c) Inference: Generate new speech with the adapted model - Anti-spoofing is an important area for continued research.



0.6

**False Positive Rate** 

0.8

0.6

0.4

0.2

True Positive Rate

**Observation:** partially linearly separable on single speaker d-vectors



Linear SVM on d-vectors (85% accuracy)

Real vs Synthesized

p334

Generated

Real

p294

**%** 

d-vectors (t-SNE)

p345

p299

p311

p305

p300



#### **Beyond the Paper: Faster Inference Models**

- It is possible to use a WaveRNN model to achieve faster inference, without compromising on sample quality.
- Few-shot adaptation for WaveRNN works out-of-the-box.



#### **Experimental Procedure**

- Train: - More than 2300+ speakers.
  - 300-hour of Google American English TTS corpus.
  - 500-hour public LibriSpeech corpus.
- Adapt on a hold-out speaker for: Adapt: - LibriSpeech: 10 sec - 5 min; VCTK: 10 sec - 10 min
- Evaluate: -(Subjective) Mean Opinion Score (MOS) human ratings:
  - for naturalness, and
  - for speaker similarity.
  - (Objective) Speaker embedding (*d*-vectors) similarity from
    - Google's Text-independent Speaker Verification model (Wan et al., 2018).



- Impressive performance even with only 10 seconds of audio from new speakers.
- State-of-the-art performance in naturalness and voice similarity.
- Applicable to restoring the voices of speech-impaired patients.
- Promising results on detection of synthesized voices.