Al Human: Large-Scale Text-to-Speech Applications

Eunwoo Song / NAVER Cloud

Who am I?

Developing Text-to-Speech engines for

- NAVER: Maps, Dictionary, News
- NAVER CLOVA: AI Speaker, Dubbing Editor
- NAVER Cloud: Papago Translator, Voice API





Deep learning-based TTS system



Human-like voice quality 🙄

Deep learning-based TTS system





Acoustic model + Vocoding model

Deep learning-based TTS system





Estimating acoustic parameters from text inputs

Speaker-specific attributes (tone, volume, timbre, speaking rate, ...)

Acoustic model + Vocoding model

Deep learning-based TTS system





Estimating speech signals from acoustic parameters

Acoustic model + Vocoding model

Deep learning-based TTS system





Acoustic parameters..?

Speaker-specific attributes (tone, volume, timbre, speaking rate, ...)

Speech production model



https://www.youtube.com/watch?v=X_JvfZiGEek

- Vocal cords
 - Voiced sound : quasi-periodic
 - Unvoiced sound : noisy
 - → 목소리의 톤을 결정 (아↘아↗)
- Vocal tract
 - Shaping voice color
 - → 발음을 결정 (아/에/이/오/우)

Speech analysis



음성 신호는 시간 축에서 특정한 에너지를 갖는 파형의 형태로 존재합니다

Speech analysis



Fourier 변환을 통해 주파수 축에서 음성을 관찰할 수 있습니다

Speech analysis



F0: 목소리의 톤을 표현하는 파라미터 (아↘아↗)

Speech analysis



Formant: 발음을 표현하는 파라미터 (아/에/이/오/우)

Speech analysis



http://hyperphysics.phy-astr.gsu.edu/hbase/Music/vowel.html

Formant: 발음을 표현하는 파라미터 (아/에/이/오/우)

Speech analysis





Mel-spectrogram: 음성의 다양한 특성들을 시간-주파수 축으로 표현

Deep learning-based TTS system





Acoustic parameters..?

Speaker-specific attributes (tone, volume, timbre, speaking rate, ...)

Deep learning-based TTS system





Estimating acoustic parameters from text inputs

Speaker-specific attributes (tone, volume, timbre, speaking rate, ...)

주어진 입력 텍스트로 부터 사람의 음성 특성을 모델링 하는 태스크

How to generate acoustic parameters?



How to generate acoustic parameters?



Text analyzer extracts phoneme sequence from the given text

음소: 음운론상의 최소 단위

- 3.2km → 삼쩌미키로미터 → ㅅ/ㅏ/ㅁ/ㅉ/ㅓ/ㅁ/ㅣ/ㅋ/ㅣ/ㄹ/ㅗ/ㅁ/ㅣ/ㅌ/ㅓ
- naver.com → 네이버닫컴 → ㄴ/ㅔ/ㅣ/ㅂ/ㅓ/ㄷ/ㅏ/ㄷ/ㅋ/ㅓ/ㅁ
- 1588-7942 → 이로팔팔 칠구사이 → |/ㄹ/ㅗ/ㅍ/ㅏ/ㄹ/ㅍ/ㅏㄹ/ㅊ/|/ㄹ/ㄱ/ㅜ/ㅅ/ㅏ/|

How to generate acoustic parameters?



Linguistic encoder extracts high-level context features from the given text

How to generate acoustic parameters?



Duration model predicts phoneme duration to upsample context embeddings into frame-level

How to generate acoustic parameters?



Duration model predicts phoneme duration to upsample context embeddings into frame-level

How to generate acoustic parameters?



Acoustic decoder predicts acoustic parameters from the given context embeddings

How to generate acoustic parameters?



Trained by **speech-text pair** recorded by single speaker or multiple speakers

Multi-speaker TTS

Speaker-dependent training



Training TTS acoustic model using target speaker's speech corpus

Speaker-dependent training



Training TTS acoustic model using target speaker's speech corpus

This approach is highly effective as the model can learn the target speaker's speaking patterns when plenty amount of recorded data is available (> 10 hours)

Speaker-dependent training



Training TTS acoustic model using target speaker's speech corpus The synthetic quality is significantly degraded when the amount of available recorded data is limited

Multi-speaker training



Training TTS acoustic model using speech corpora from multiple speakers

Multi-speaker training



Training TTS acoustic model using speech corpora from multiple speakers Some layers capture speaker-independent characteristics shared across different speakers Other layers represent speaker-dependent characteristics specific to the target speaker

Multi-speaker training



Training TTS acoustic model using speech corpora from multiple speakers

This offers a promising solution to address the data shortage problem in speaker-dependent TTS model



Training TTS acoustic model using speech corpora from multiple speakers with speaker embeddings

Linguistic encoder capture speaker-independent contents shared across different speakers



Training TTS acoustic model using speech corpora from multiple speakers with speaker embeddings

Linguistic encoder capture speaker-independent contents shared across different speakers

Acoustic decoder and duration model represent speaker-dependent timbre and prosody, respectively



Multi-speaker model



How to design the speaker embeddings ?



Speaker recognition model is designed to capture the speaker identity from the given acoustic parameter

Its hidden representations contain speaker-specific characteristics



Multi-speaker model


Multi-style TTS

Style-dependent training



Training TTS acoustic model using speech corpus containing target speaking style This approach is highly effective as the model can learn the target speaking style when plenty amount of recorded data is available (> 10 hours)

Style-dependent training



Training TTS acoustic model using speech corpus containing target speaking style

The synthetic quality is significantly degraded

when the amount of available recorded data is limited

Multi-style training



Training TTS acoustic model using speech corpora from multiple styles Linguistic encoders capture style-independent characteristics shared across different styles Acoustic decoder and duration model represent style-dependent characteristics







VAE-based style encoder extracts style representations from the given acoustic parameters



Multi-style model

VAE-based style encoder





Multi-style model

VAE-based style encoder μ Z Style Acoustic -> -> Encoder Decoder σ^2 Mean vector μ Unsupervised style representation σ^2 Variance vector Latent vector ~ $N(\mu, \sigma^2)$ Ζ

t-SNE plot of VAE representations (= μ) obtained from multiple emotions



Multi-style model



t-SNE plot of VAE representations (= μ) obtained from multiple emotions



Multi-style model



t-SNE plot of VAE representations (= μ) obtained from multiple emotions



Multi-style model



Acoustic model

Large-scale TTS

Recording constraint

	Conventional TTS	Custom TTS
Speaker	Professional	Non-professional
Environment	Clean studio	Anywhere
Amount	> 30~60 min	< Few seconds
Speaking type	Script reading	Spontaneous
Synthetic quality	Very natural	Unnatural

Recording constraint

	Conventional TTS	Custom TTS
Speaker	Professional	Non-professional
Environment	Clean studio	Anywhere
Amount	> 30~60 min	< Few seconds
Speaking type	Script reading	Spontaneous
Synthetic quality	Very natural	Natural

The trend is shifting towards scaling up TTS model using large-scale dataset

Voicebox model from Meta

Voicebox: Text-Guided Multilingual Universal Speech Generation at Scale

Matthew Le* Apoorv Vyas* Bowen Shi* Brian Karrer* Leda Sari Rashel Moritz Mary Williamson Vimal Manohar Yossi Adi[†] Jay Mahadeokar Wei-Ning Hsu*

Fundamental AI Research (FAIR), Meta

Abstract

Large-scale generative models such as GPT and DALL-E have revolutionized natural language processing and computer vision research. These models not only generate high fidelity text or image outputs, but are also generalists which can solve tasks not explicitly taught. In contrast, speech generative models are still primitive in terms of scale and task generalization. In this paper, we present Voicebox, the most versatile text-guided generative model for speech at scale. Voicebox is a non-autoregressive flow-matching model trained to infill speech, given audio context and text, trained on over 50K hours of speech that are neither filtered nor enhanced. Similar to GPT, Voicebox can perform many different tasks through in-context learning, but is more flexible as it can also condition on future context. Voicebox can be used for mono or cross-lingual zero-shot text-to-speech synthesis, noise removal, content editing, style conversion, and diverse sample generation. In particular, Voicebox outperforms the state-of-the-art zero-shot TTS model VALL-E on both intelligibility (5.9% vs 1.9% word error rates) and audio similarity (0.580 vs 0.681) while being up to 20 times faster. Audio samples can be found in https://voicebox.metademolab.com.

Self-supervised learning representation

Model learns useful representations of data without relying on manually labeled data



Self-supervised learning representation

Model learns useful representations of data without relying on manually labeled data





Self-supervised learning representation

Model learns useful representations of data without relying on manually labeled data



Self-supervised learning representation

Model learns useful representations of data without relying on manually labeled data



Self-supervised learning representation

Model learns useful representations of data without relying on manually labeled data



Self-supervised learning representation

Model learns useful representations of data without relying on manually labeled data











Voice cloning



Voice cloning



Voice cloning



디코더는 음성 프롬프트의 특성을 반영하여 TTS 생성을 할 수 있음

Voice cloning



디코더는 음성 프롬프트의 특성을 반영하여 TTS 생성을 할 수 있음

Voice cloning

	Conventional TTS	Custom TTS
Speaker	Professional	Non-professional
Environment	Clean studio	Anywhere
Amount	> 30~60 min	< Few seconds
Speaking type	Script reading	Spontaneous
Synthetic quality	Very natural	Natural <u>Unnatural</u>

Voice cloning



Summary





Human-like voice quality 🙄







Acoustic model + Vocoding model





Estimating acoustic parameters from text inputs

Speaker-specific attributes (tone, volume, timbre, speaking rate, ...)

Acoustic model + Vocoding model

Acoustic model + Vocoding model

Summary

Data collection

	Conventional TTS	Custom TTS
Speaker	Professional	Non-professiona
Environment	Clean studio	Anywhere
Amount	> 30~60 min	< Few seconds
Speaking type	Script reading	Spontaneous
Synthetic quality	Very natural	Natural

Multi-speaker & Multi-style TTS

Large-scale TTS





gregorio.song@gmail.com

