# CUAI 딥러닝 논문 리뷰 스터디 (음성 처리)
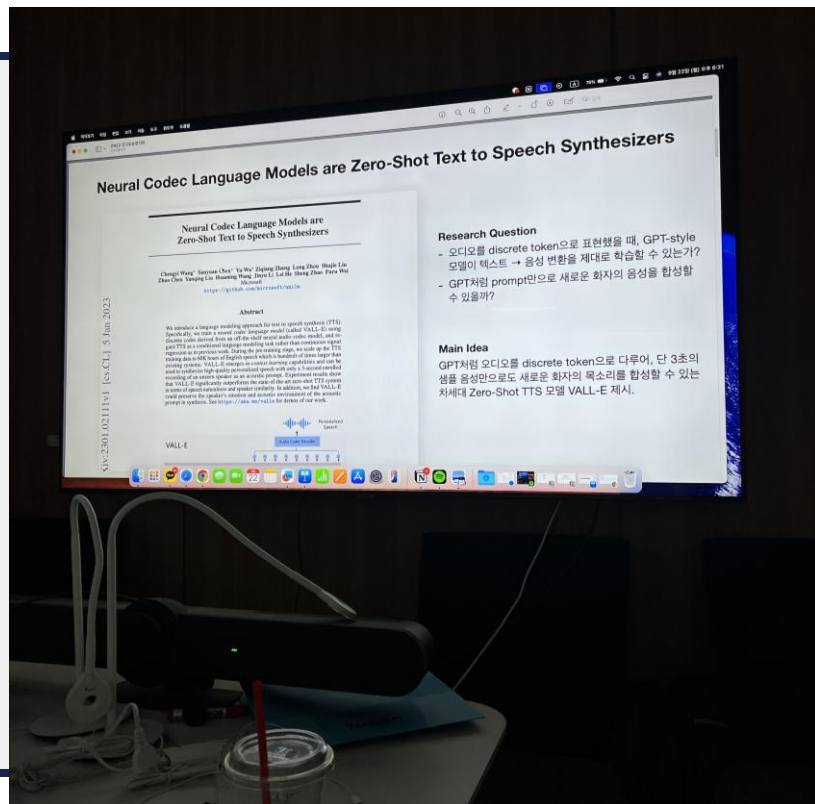
2025.09.30

발표자 : 김동건

# 스터디원 소개 및 만남 인증



스터디원 1 : 김동건

스터디원 2 : 양희원

스터디원 3 : 이나현

# 스터디 방식

**월요일 15:00 – 16:30**
**대면 스터디 진행**


**각자 공부할 논문 선정 후 간단하게 PPT 제작 후**
**스터디에서 발표하며 공부한 논문의 내용을 공유**

# 스터디 주제

**9월 22일 월요일 15:00 – 16:30 1차 스터디 진행함**
**각자 공부할 논문 선정 후 간단한 설명 완료.**
**10월 13일까지 공부 완료하여 10월 13일 2차 스터디에 공부 내용 공유할 예정.**

**공부할 논문**

- **김동건:** ProMode: A Speech Prosody Model Conditioned on Acoustic and Textual Inputs
- **양희원:** Long-Form Speech Generation with Spoken Language Model
- **이나현:** Neural Codec Language Models are Zero-Shot Text to Speech Synthesizers

# 스터디 내용



TTS 합성 시 음성의 자연스러움을 결정하는 Prosody에 대한 임베딩을 미리 학습하여, SER, TTS 등 Downstream Task의 종류에 상관없이 Prosody의 정확도를 높이는 모듈에 대한 내용

내용

# 스터디 내용



**Long-Form Speech Generation with Spoken Language Models**

Se Jin Park [* 1 2]  Julian Salazar [* 1]  Aren Jansen [1]  Keisuke Kinoshita [1]  Yong Man Ro [2]  RJ Skerry-Ryan [1]

## Abstract

We consider the generative modeling of speech over multiple minutes, a requirement for long-form multimedia generation and audio-native voice assistants. However, textless spoken language models struggle to generate plausible speech past tens of seconds, due to high temporal resolution of speech tokens causing loss of coherence, architectural issues with long-sequence training or extrapolation, and memory costs at inference time. From these considerations we derive **SpeechSSM**, the first speech language model family to learn from and sample long-form spoken audio (e.g., 16 minutes of read or extemporaneous speech) in a single decoding session without text intermediates. SpeechSSMs leverage recent advances in linear-time sequence modeling to greatly surpass current Transformer spoken LMs in coherence and efficiency on multi-minute generations while still matching them at the utterance level. As we found current spoken language evaluations uninformative, especially in this new long-form setting, we also introduce: **LibriSpeech-Long**, a benchmark for long-form speech evaluation; new embedding-based and LLM-judged metrics; and quality measurements over length and time. Speech samples, the LibriSpeech-Long dataset, and any future code or model releases can be found at https://google.github.io/tacotron/publications/speechssm/.
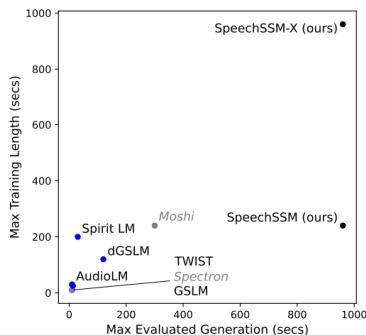
*Figure 1.* Maximum sequence lengths considered by various spoken LMs. *Italicized* models used text intermediates at generation time. Our models can generate indefinitely due to their constant memory footprint, but we cap our evaluations to 16 minutes.

## 1. Introduction

its paralinguistic aspects, such as prosody (Kharitonov et al., 2022) and turn-taking (Nguyen et al., 2023b). These capabilities make speech-native language models (LMs) promising for applications like media understanding and co-creation, audio-native voice assistants, and textless NLP. However, real-world use-cases of spoken LMs require the ability to both understand and generate long-form speech. For example, voice interactions can last many minutes, requiring a model to maintain a growing conversational history in real time, and expressive media like audiobooks and podcasts can require semantic, paralinguistic, and speaker coherence

**ICML 2025 Oral**

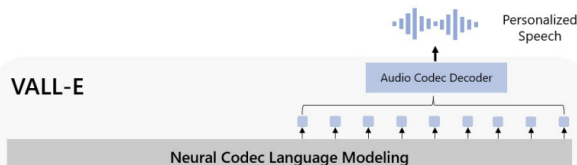수십분의 긴 음성 생성을 목표로, SSM을 활용해 텍스트를 거치지 않고 직접 음성을 생성하는 모델을 제안함.

# 스터디 내용

## Neural Codec Language Models are Zero-Shot Text to Speech Synthesizers

Chengyi Wang* Sanyuan Chen* Yu Wu* Ziqiang Zhang Long Zhou Shujie Liu
Zhuo Chen Yanqing Liu Huaming Wang Jinyu Li Lei He Sheng Zhao Furu Wei
Microsoft
https://github.com/microsoft/unilm

### Abstract

We introduce a language modeling approach for text to speech synthesis (TTS). Specifically, we train a *neural codec language model* (called VALL-E) using discrete codes derived from an off-the-shelf neural audio codec model, and regard TTS as a conditional language modeling task rather than continuous signal regression as in previous work. During the pre-training stage, we scale up the TTS training data to 60K hours of English speech which is hundreds of times larger than existing systems. VALL-E emerges *in-context learning* capabilities and can be used to synthesize high-quality personalized speech with only a 3-second enrolled recording of an unseen speaker as an acoustic prompt. Experiment results show that VALL-E significantly outperforms the state-of-the-art zero-shot TTS system in terms of speech naturalness and speaker similarity. In addition, we find VALL-E could preserve the speaker's emotion and acoustic environment of the acoustic prompt in synthesis. See https://aka.ms/valle for demos of our work.

**Microsoft 2023**

텍스트를 음성 코드로 바로 변환해 Zero-Shot TTS 를 가능하게 하고, 짧은 음성 샘플만으로 화자의 음색을 모방할 수 있고, 별도 학습 없이도 새로운 화자와 문장을 자연스럽게 합성할 수 있는 VALL-E를 제안함