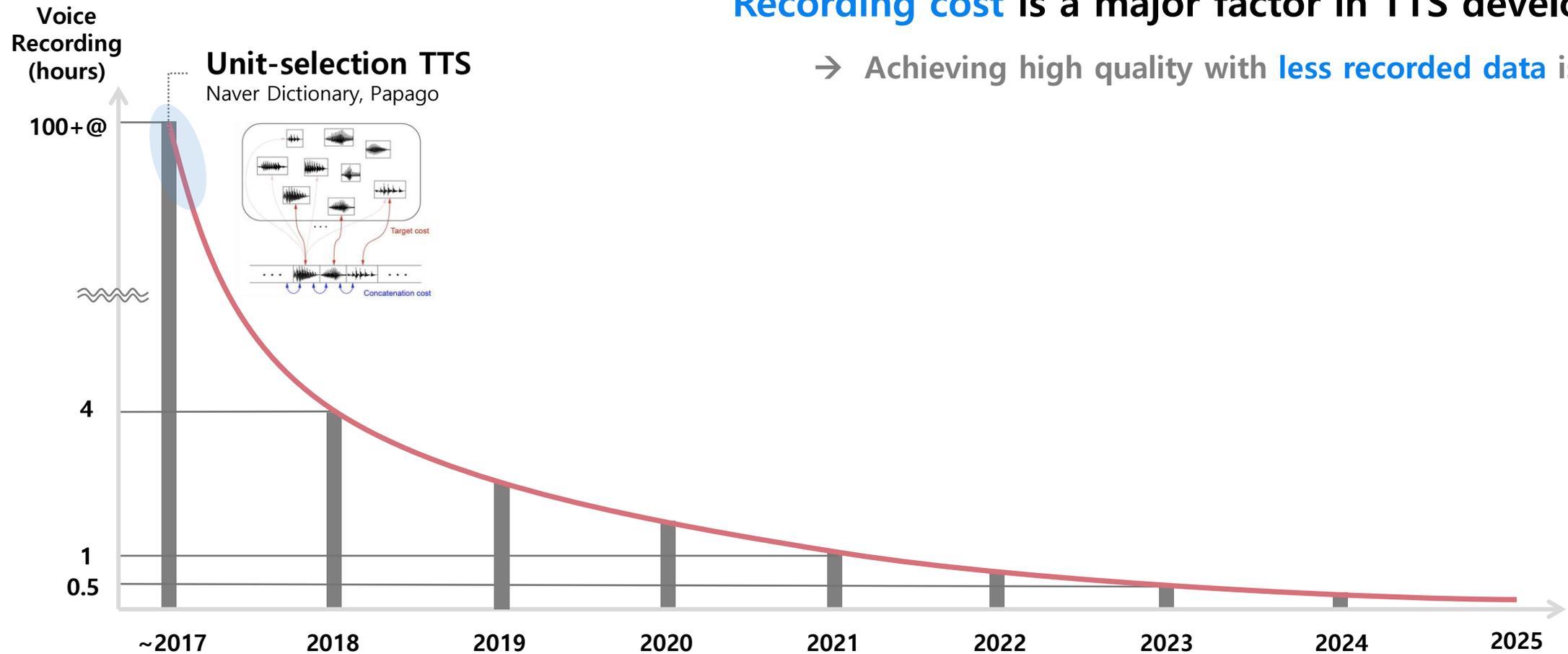


# Zero-shot Voice Cloning

Eunwoo Song / NAVER Cloud

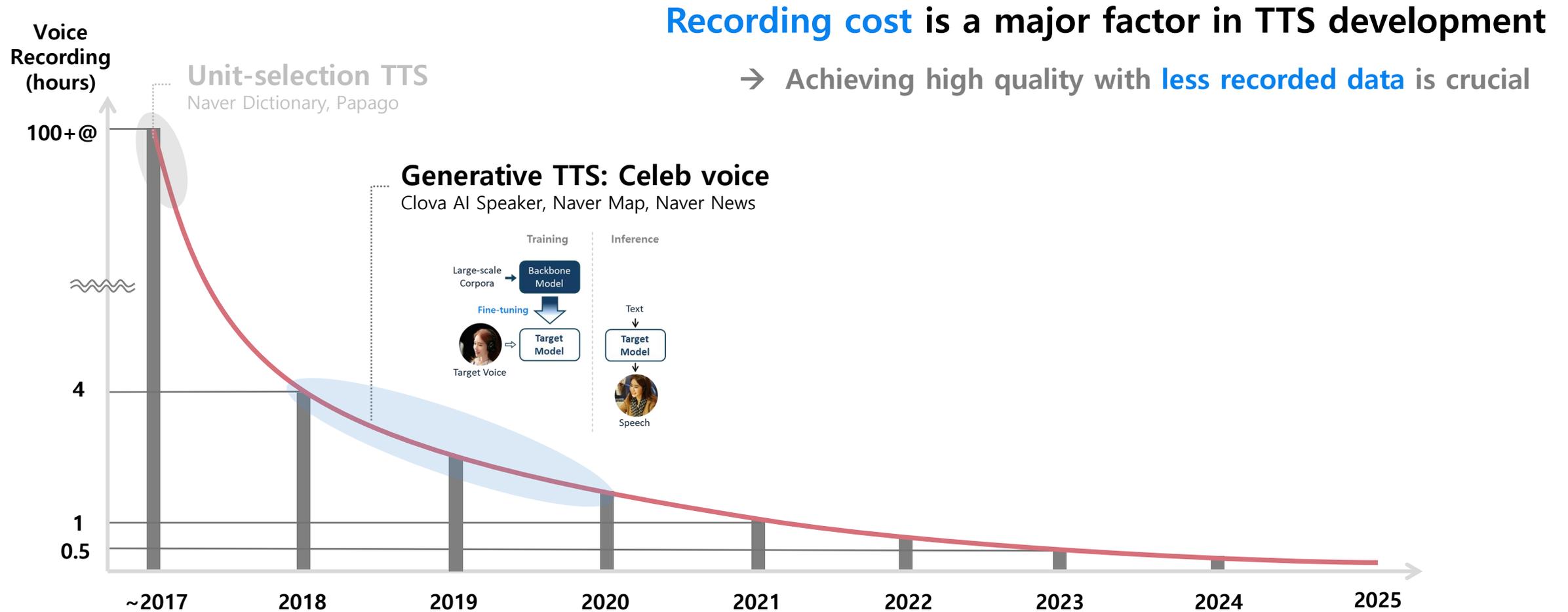
# Milestone



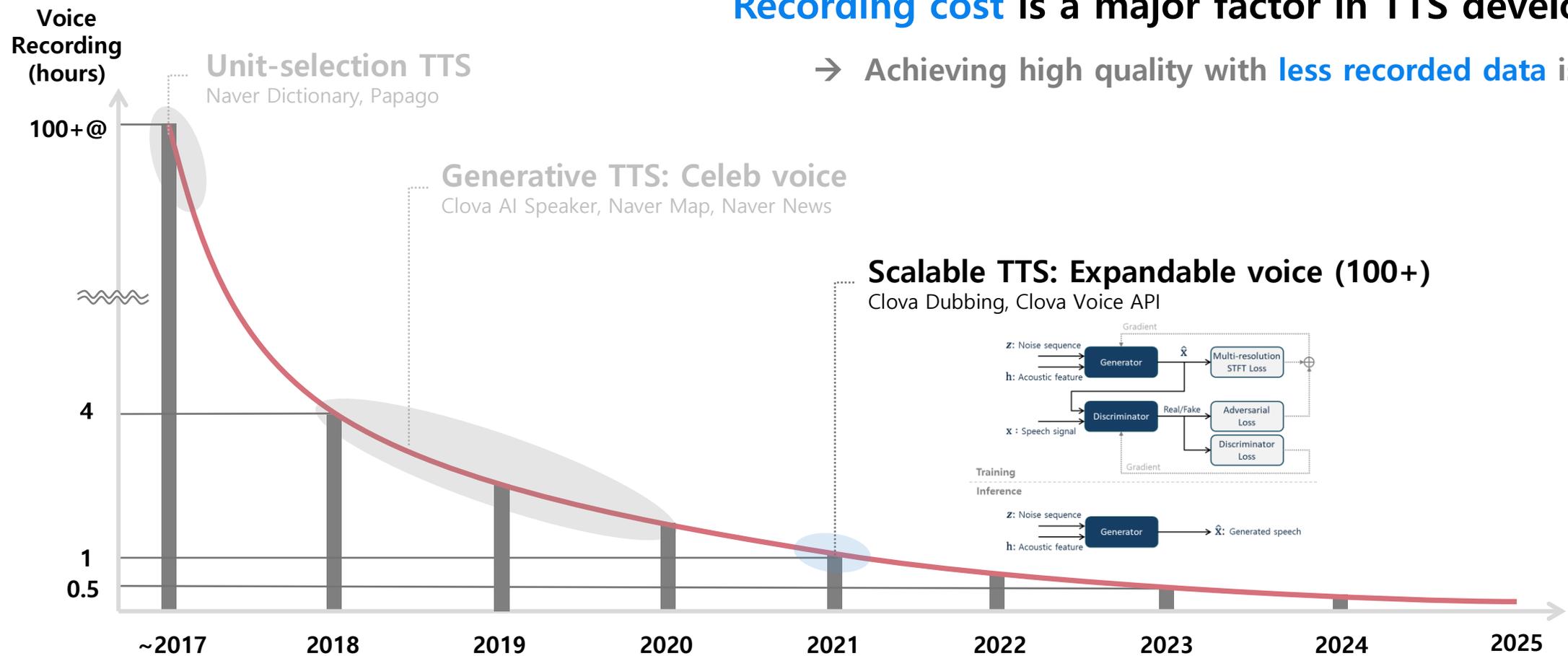
**Recording cost** is a major factor in TTS development

→ Achieving high quality with **less recorded data** is crucial

# Milestone

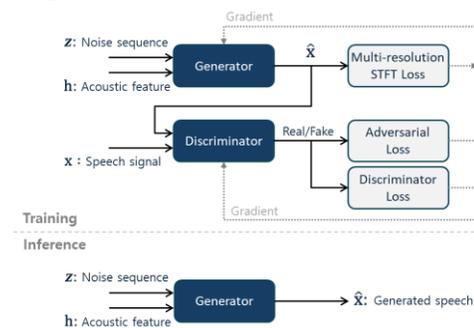


# Milestone

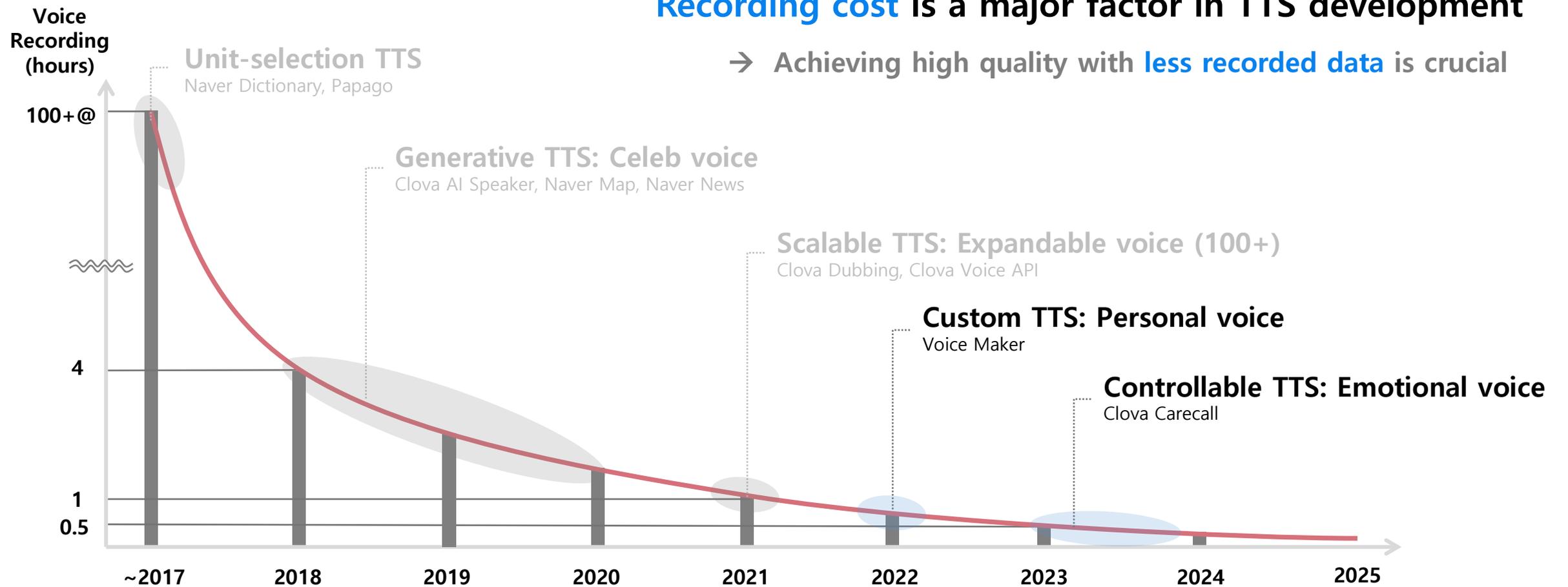


**Recording cost** is a major factor in TTS development

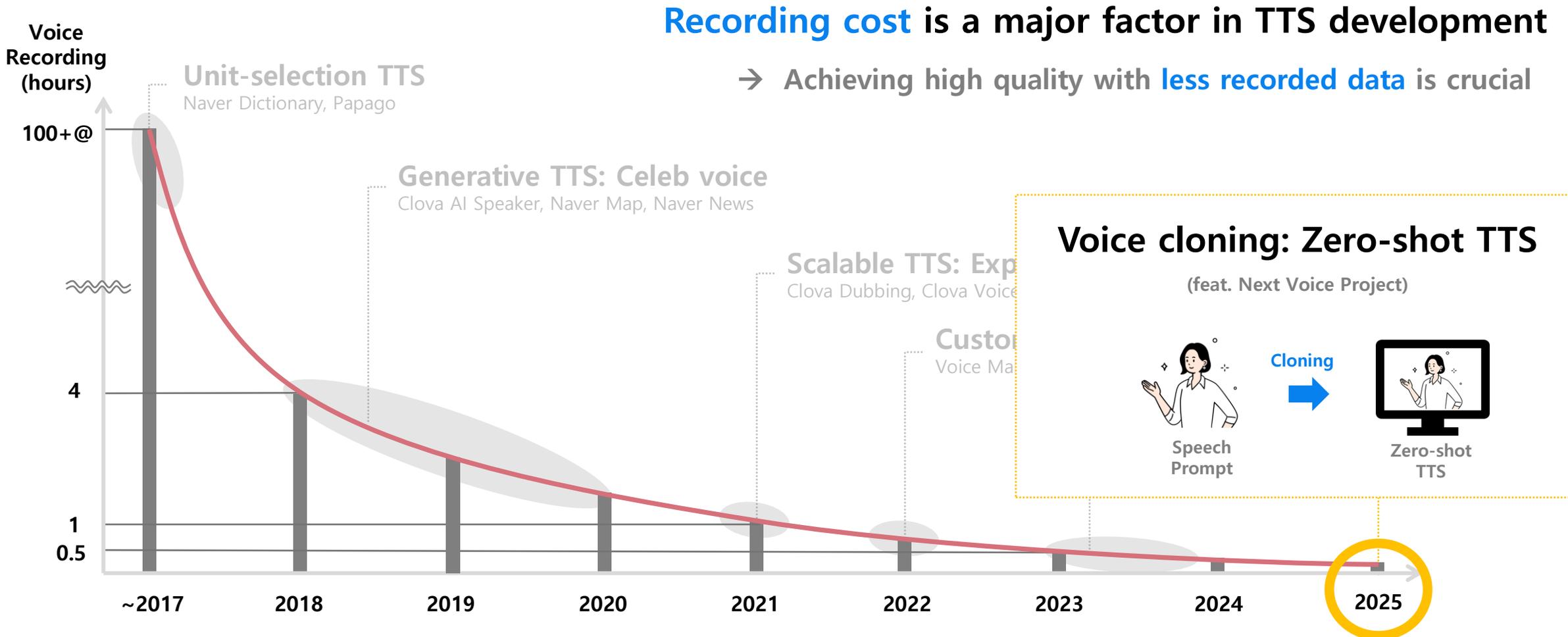
→ Achieving high quality with **less recorded data** is crucial



# Milestone



# Milestone



# Zero-shot Voice Cloning

---

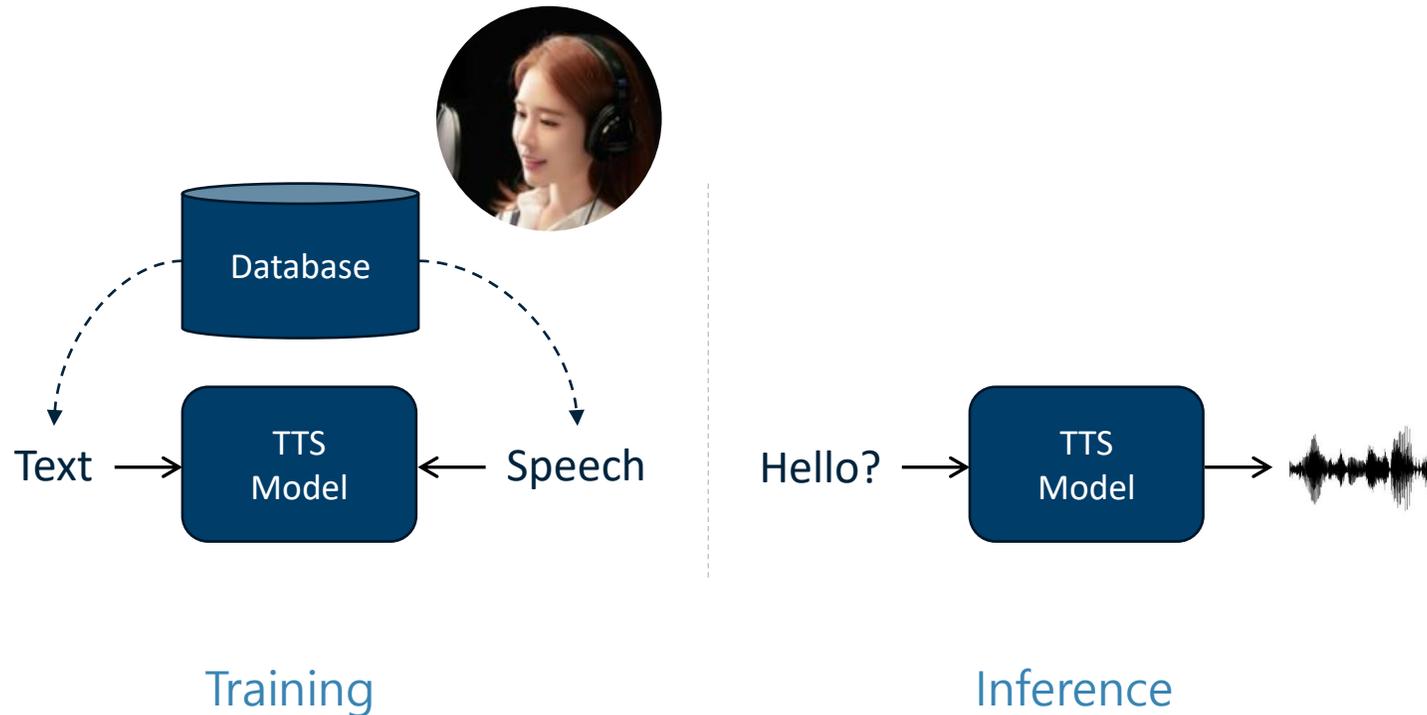
# Contents

---

1. Introduction
2. Speech analysis method
3. Text-to-speech acoustic model
4. Zero-shot voice cloning

# Introduction

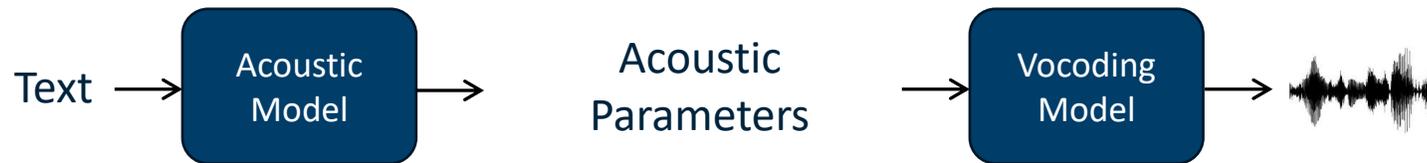
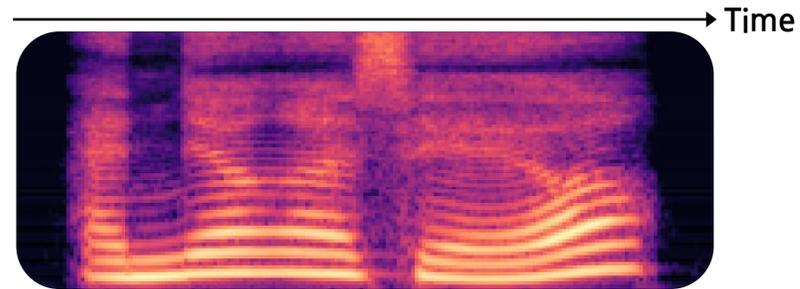
## Deep learning-based TTS system



Human-like voice quality 😊

# Introduction

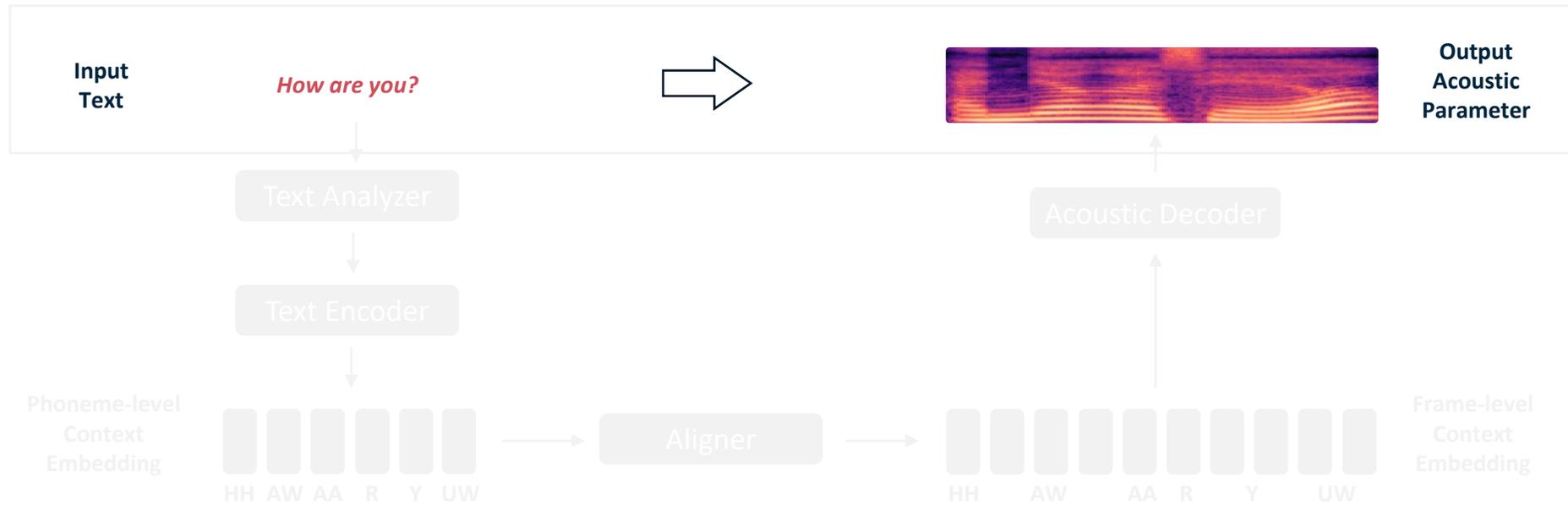
Deep learning-based TTS system



**Acoustic model + Vocoding model**

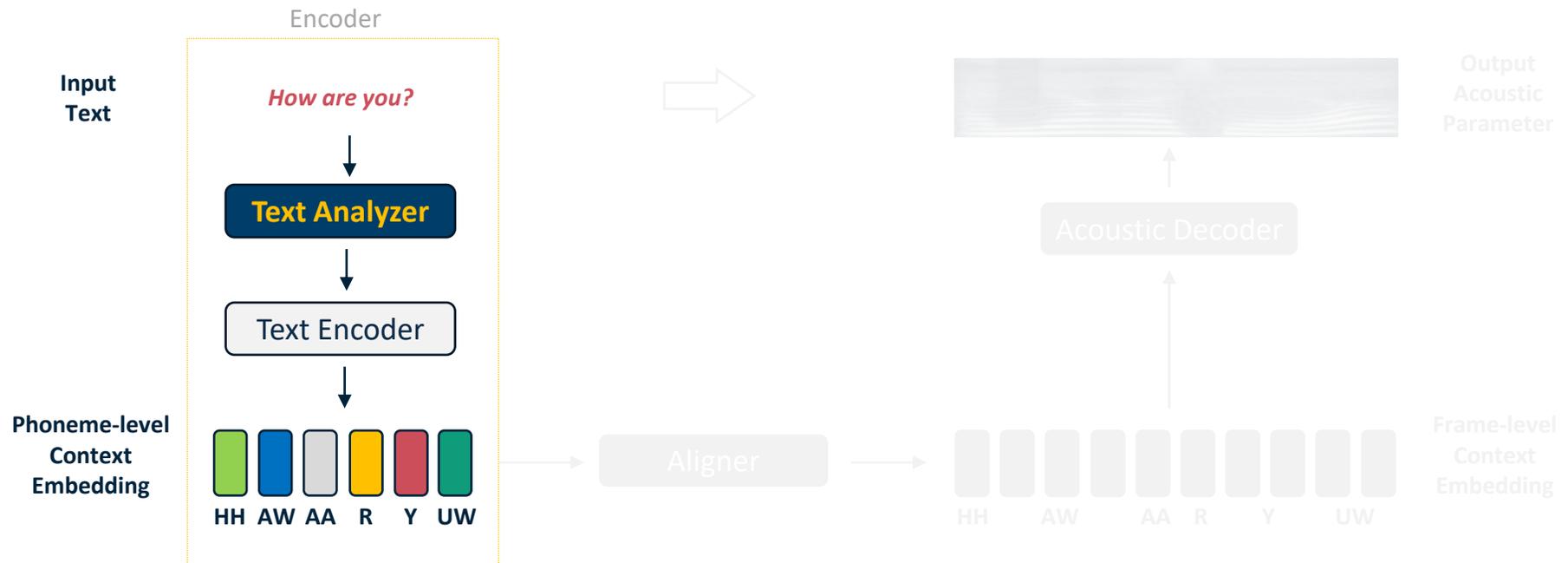
# TTS acoustic model

How to generate acoustic parameters?



# TTS acoustic model

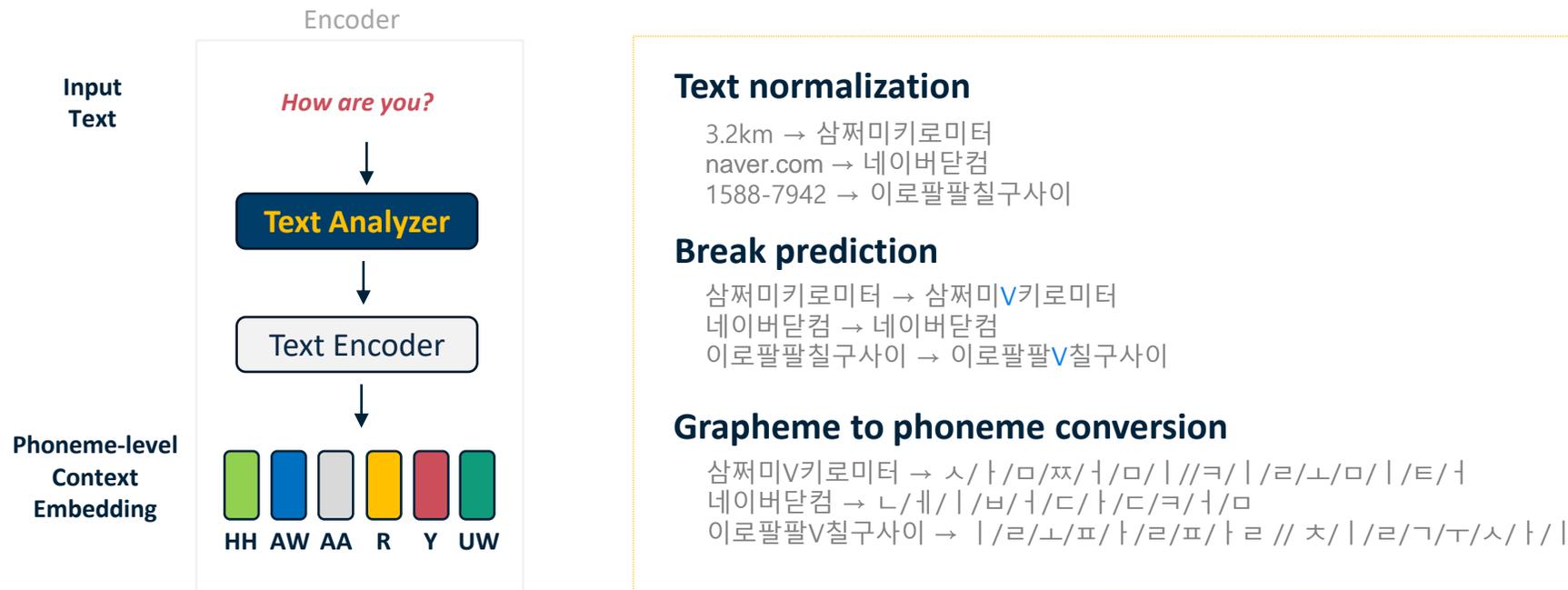
How to generate acoustic parameters?



**Text analyzer** extracts **phoneme** sequence from the given text

# TTS acoustic model

## How to generate acoustic parameters?

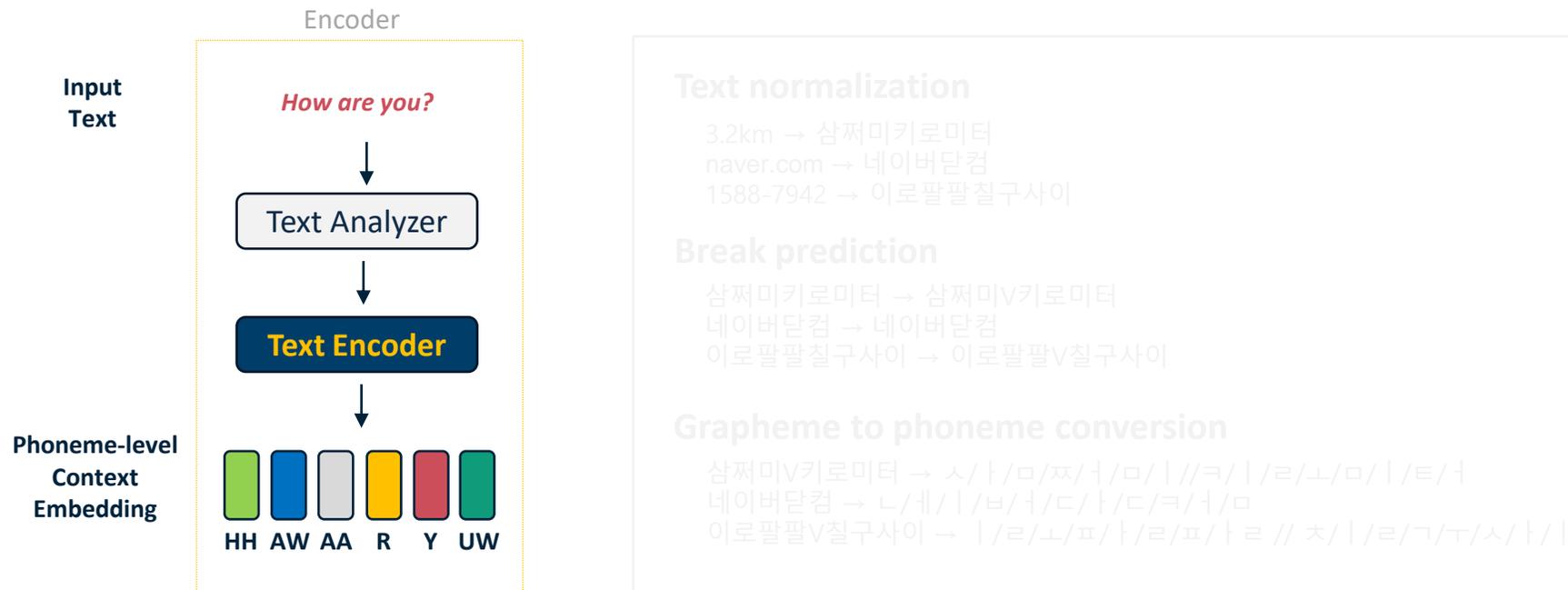


Text analyzer extracts **phoneme** sequence from the given text

음소: 음운론상의 최소 단위

# TTS acoustic model

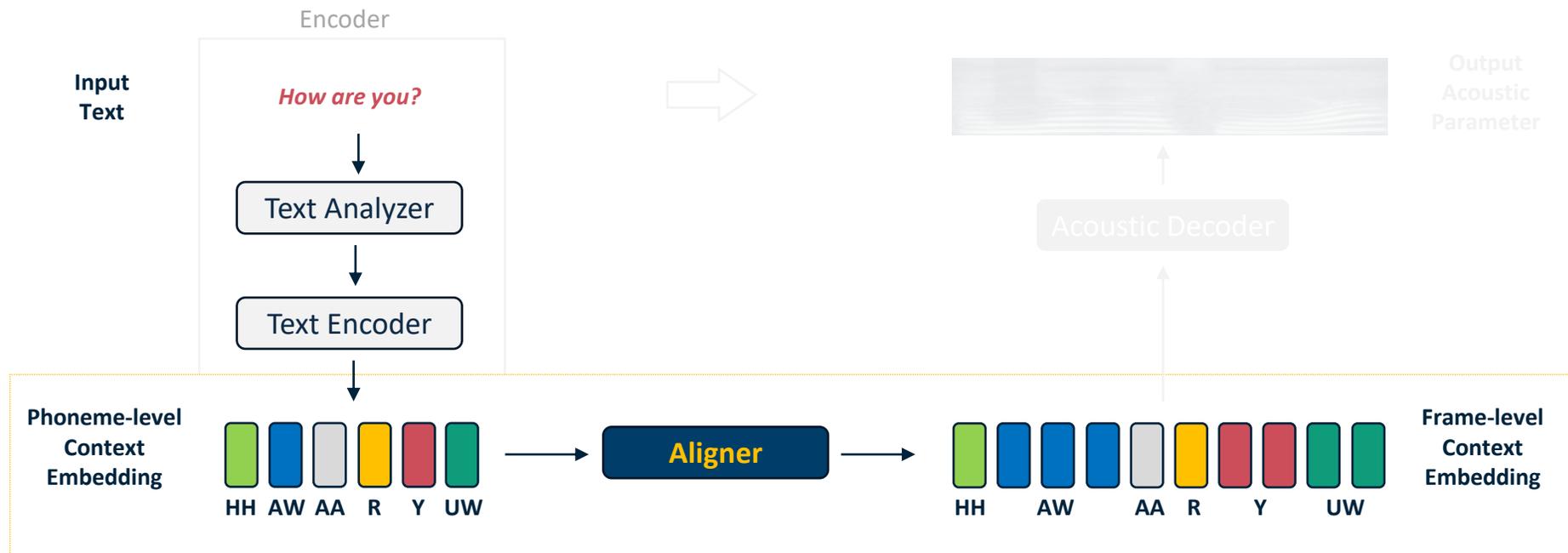
How to generate acoustic parameters?



**Text encoder** extracts high-level **context features** from the given phoneme sequence

# TTS acoustic model

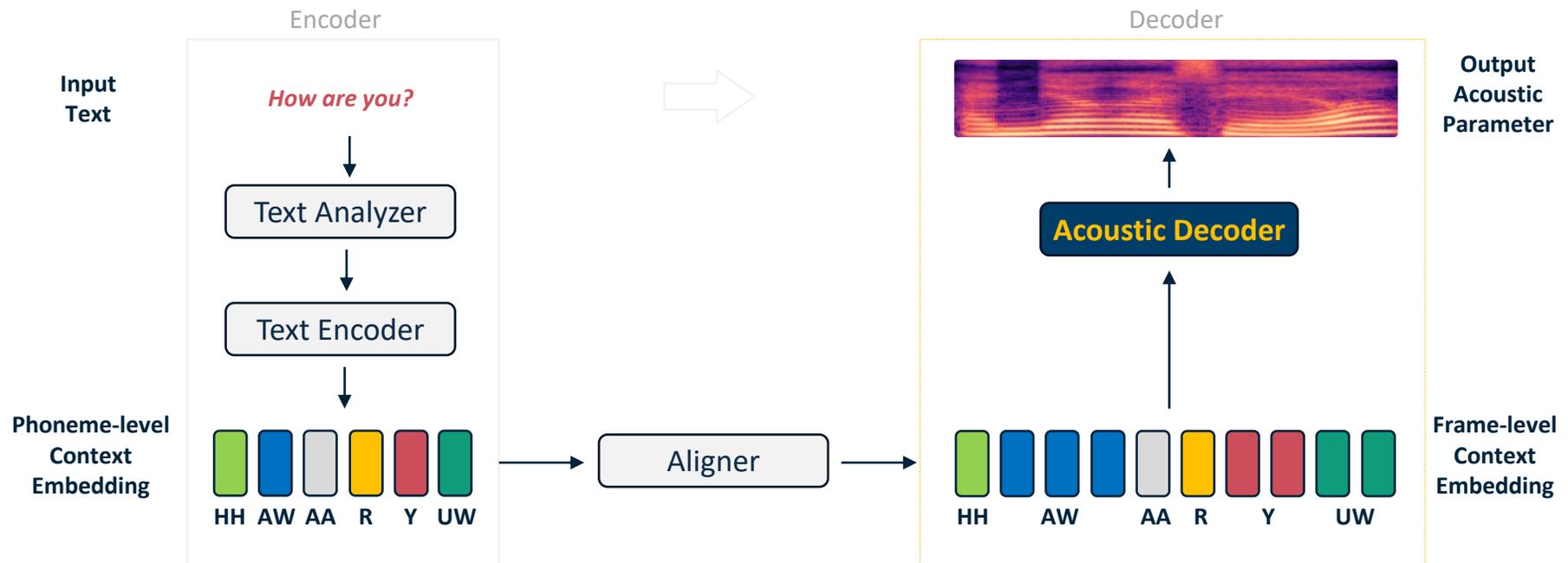
How to generate acoustic parameters?



**Aligner** upamples context embeddings from **phoneme-level** to **frame-level**

# TTS acoustic model

How to generate acoustic parameters?



**Acoustic decoder** predicts **acoustic parameters** from the given context embeddings

# TTS acoustic model

How to generate acoustic parameters?



Acoustic decoder predicts acoustic parameters from the given context embeddings

# TTS acoustic model

Statistical parametric speech synthesis (2013)

**STATISTICAL PARAMETRIC SPEECH SYNTHESIS USING DEEP NEURAL NETWORKS**

*Heiga Zen, Andrew Senior, Mike Schuster*

Google

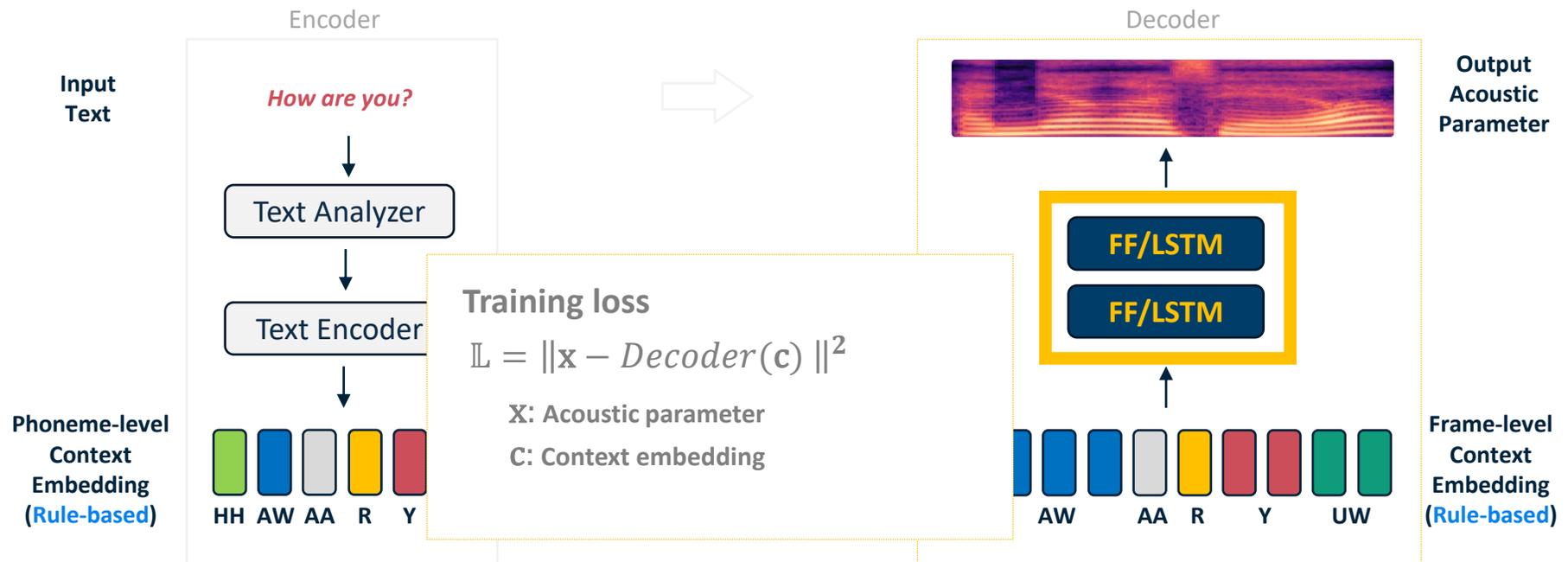
{heigazen, andrewsenior, schuster}@google.com

## **ABSTRACT**

Conventional approaches to statistical parametric speech synthesis typically use decision tree-clustered context-dependent hidden Markov models (HMMs) to represent probability densities of speech parameters given texts. Speech parameters are generated from the probability densities to maximize their output probabilities, then a speech waveform is reconstructed from the generated parameters. This approach is reasonably effective but has a couple of limitations, *e.g.* decision trees are inefficient to model complex context dependencies. This paper examines an alternative scheme that is based on a deep neural network (DNN). The relationship between input texts and their acoustic realizations is modeled by a DNN. The use of the DNN can address some limitations of the conventional approach. Experimental results show that the DNN-based systems outperformed the HMM-based systems with similar numbers of parameters.

# TTS acoustic model

Statistical parametric speech synthesis (2013)



The first **DNN model** for the TTS acoustic model

# TTS acoustic model

Tacotron 2 (2018)

## NATURAL TTS SYNTHESIS BY CONDITIONING WAVENET ON MEL SPECTROGRAM PREDICTIONS

*Jonathan Shen<sup>1</sup>, Ruoming Pang<sup>1</sup>, Ron J. Weiss<sup>1</sup>, Mike Schuster<sup>1</sup>, Navdeep Jaitly<sup>1</sup>, Zongheng Yang<sup>\*2</sup>, Zhifeng Chen<sup>1</sup>, Yu Zhang<sup>1</sup>, Yuxuan Wang<sup>1</sup>, RJ Skerry-Ryan<sup>1</sup>, Rif A. Saurous<sup>1</sup>, Yannis Agiomyriannakis<sup>1</sup>, and Yonghui Wu<sup>1</sup>*

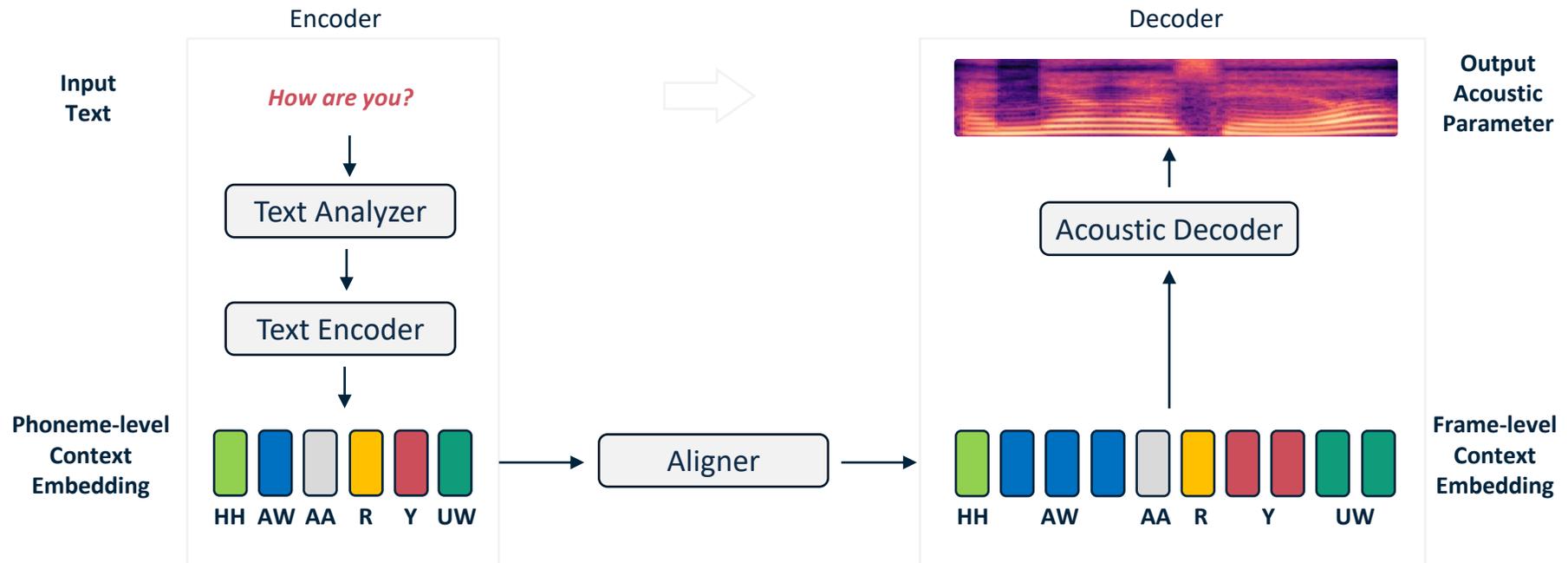
<sup>1</sup>Google, Inc., <sup>2</sup>University of California, Berkeley,  
{jonathanasdf, rpang, yonghui}@google.com

### ABSTRACT

This paper describes Tacotron 2, a neural network architecture for speech synthesis directly from text. The system is composed of a recurrent sequence-to-sequence feature prediction network that maps character embeddings to mel-scale spectrograms, followed by a modified WaveNet model acting as a vocoder to synthesize time-domain waveforms from those spectrograms. Our model achieves a mean opinion score (MOS) of 4.53 comparable to a MOS of 4.58 for professionally recorded speech. To validate our design choices, we present ablation studies of key components of our system and evaluate the impact of using mel spectrograms as the conditioning input to WaveNet instead of linguistic, duration, and  $F_0$  features. We further show that using this compact acoustic intermediate representation allows for a significant reduction in the size of the WaveNet architecture.

# TTS acoustic model

Tacotron 2 (2018)



The first **seq2seq model** for TTS acoustic model

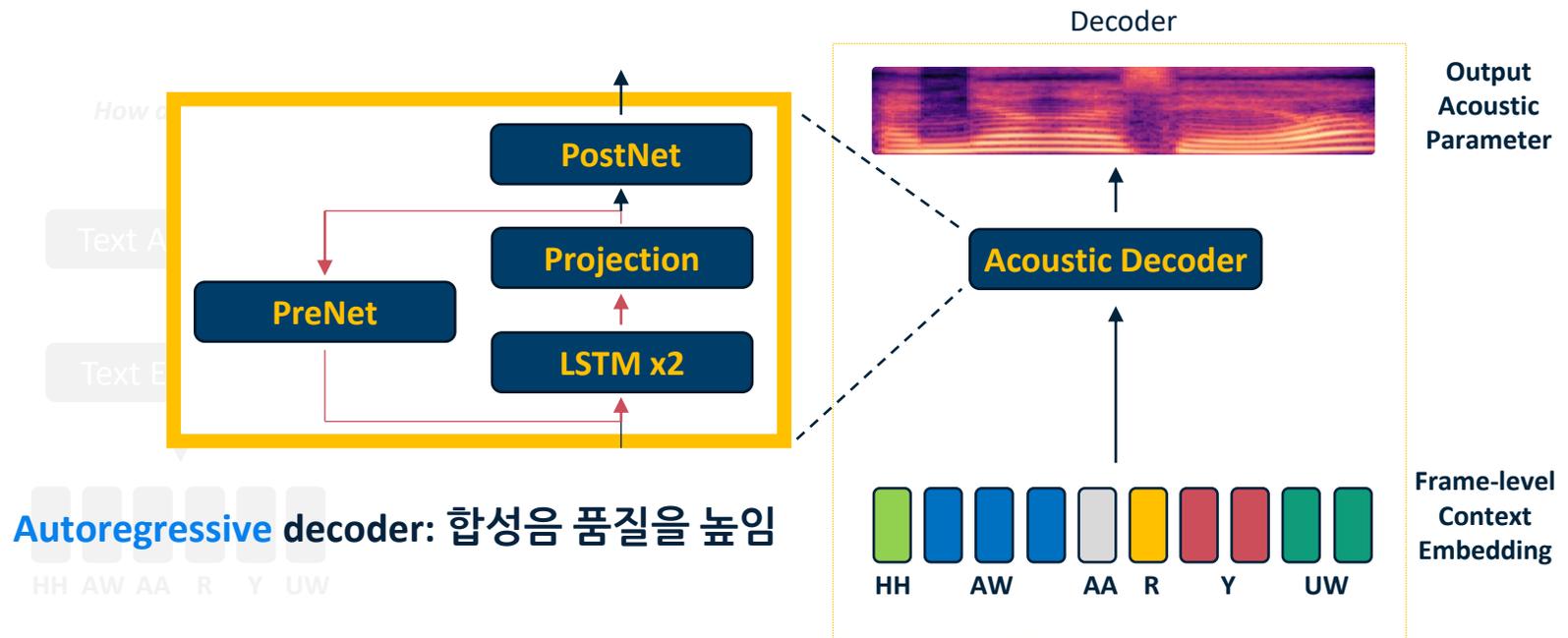
# TTS acoustic model

Tacotron 2 (2018)



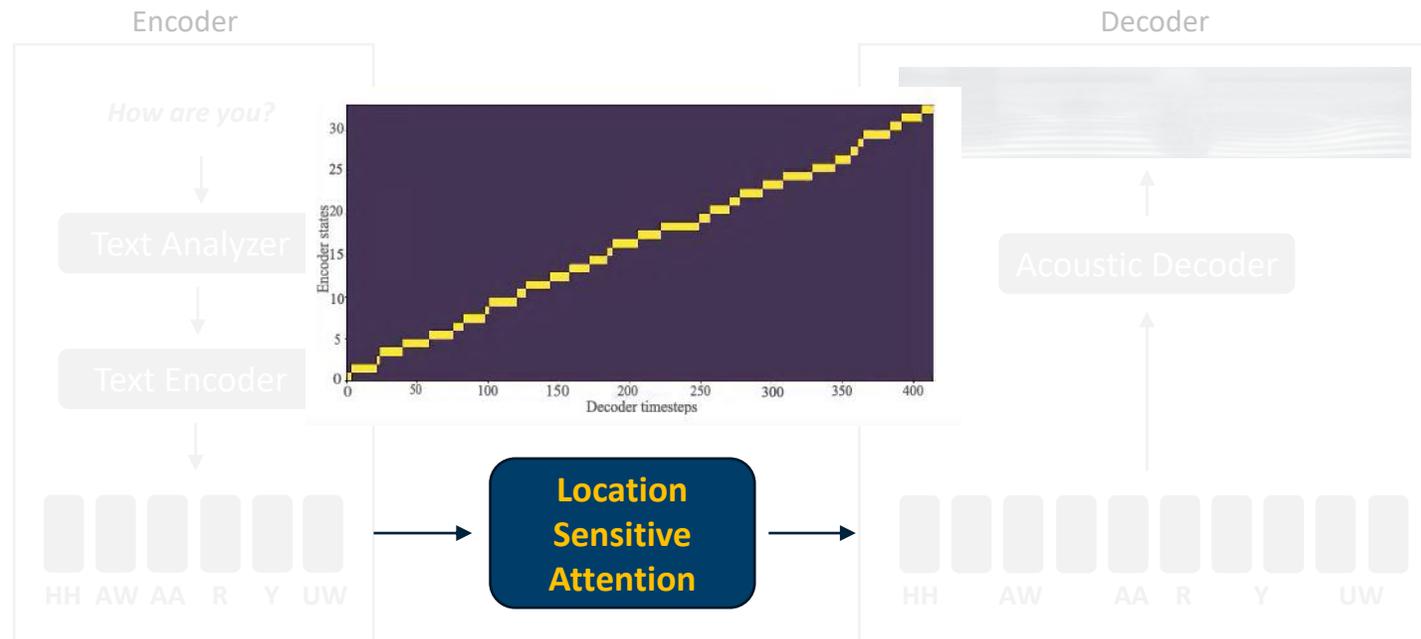
# TTS acoustic model

Tacotron 2 (2018)



# TTS acoustic model

Tacotron 2 (2018)



Attention 메커니즘을 이용해 인코더-디코더 사이의 alignment 를 얻어낼 수 있음

# TTS acoustic model

## Tacotron 2 (2018)

System	MOS
Parametric	3.492 ± 0.096
Tacotron (Griffin-Lim)	4.001 ± 0.087
Concatenative	4.166 ± 0.091
WaveNet (Linguistic)	4.341 ± 0.051
Ground truth	4.582 ± 0.053
Tacotron 2 (this paper)	<b>4.526 ± 0.066</b>

End-to-end acoustic model + WaveNet vocoder

당시 최고 합성 모델인 Concatenative 보다 우수한, 녹음에 가까운 수준의 음성 합성 모델

<https://ai.googleblog.com/2017/12/tacotron-2-generating-human-like-speech.html>

# TTS acoustic model

FastSpeech 2 (2020)

## FASTSPEECH 2: FAST AND HIGH-QUALITY END-TO-END TEXT TO SPEECH

Yi Ren<sup>1\*</sup>, Chenxu Hu<sup>1\*</sup>, Xu Tan<sup>2</sup>, Tao Qin<sup>2</sup>, Sheng Zhao<sup>3</sup>, Zhou Zhao<sup>1†</sup>, Tie-Yan Liu<sup>2</sup>

<sup>1</sup>Zhejiang University  
{rayeren, chenxuhu, zhaozhou}@zju.edu.cn

<sup>2</sup>Microsoft Research Asia  
{xuta, taoqin, tyliu}@microsoft.com

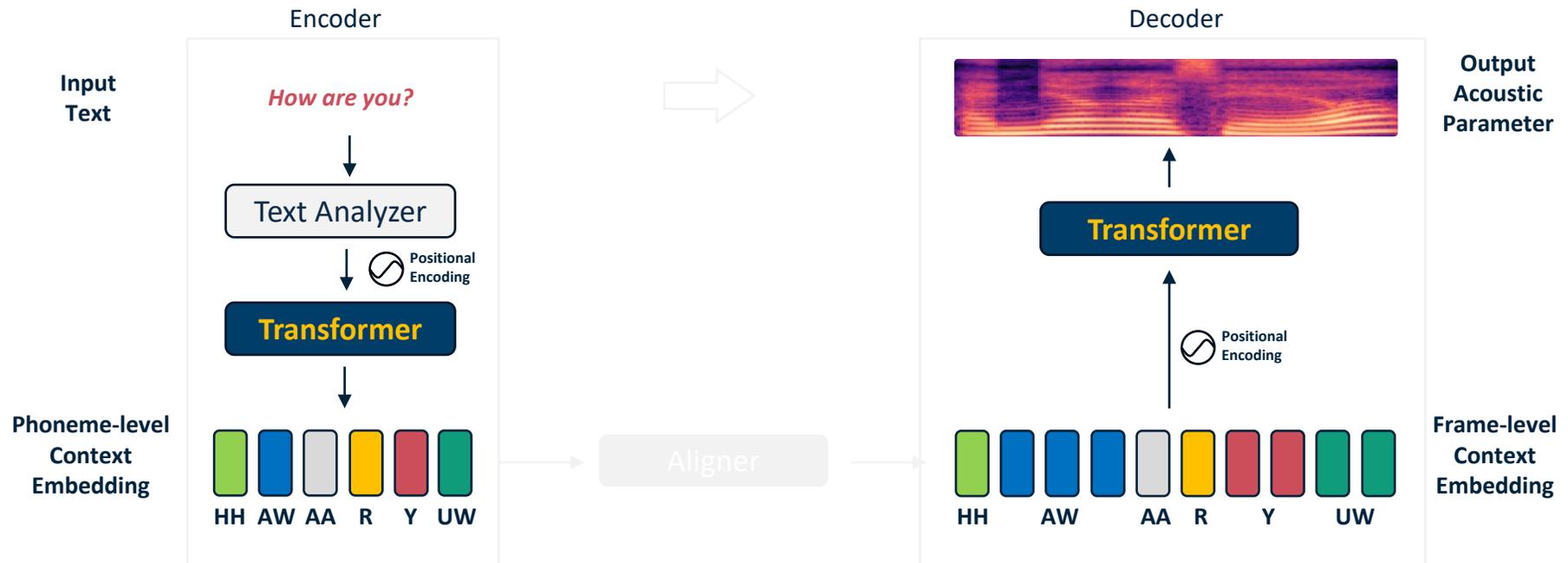
<sup>3</sup>Microsoft Azure Speech  
Sheng.Zhao@microsoft.com

## ABSTRACT

Non-autoregressive text to speech (TTS) models such as FastSpeech (Ren et al., 2019) can synthesize speech significantly faster than previous autoregressive models with comparable quality. The training of FastSpeech model relies on an autoregressive teacher model for duration prediction (to provide more information as input) and knowledge distillation (to simplify the data distribution in output), which can ease the one-to-many mapping problem (i.e., multiple speech variations correspond to the same text) in TTS. However, FastSpeech has several disadvantages: 1) the teacher-student distillation pipeline is complicated and time-consuming, 2) the duration extracted from the teacher model is not accurate enough, and the target mel-spectrograms distilled from teacher model suffer from information loss due to data simplification, both of which limit the voice quality. In this paper, we propose FastSpeech 2, which addresses the issues in FastSpeech and better solves the one-to-many mapping problem in TTS by 1) directly training the model with ground-truth target instead of the simplified output from teacher, and 2) introducing more variation information of speech (e.g., pitch, energy and more accurate duration) as conditional inputs. Specifically, we extract duration, pitch and energy from speech waveform and directly take them as conditional inputs in training and use predicted values in inference. We further design FastSpeech 2s, which is the first attempt to directly generate speech waveform from text in parallel, enjoying the benefit of fully end-to-end inference. Experimental results show that 1) FastSpeech 2 achieves a 3x training speed-up over FastSpeech, and FastSpeech 2s enjoys even faster inference speed; 2) FastSpeech 2 and 2s outperform FastSpeech in voice quality, and FastSpeech 2 can even surpass autoregressive models. Audio samples are available at <https://speechresearch.github.io/fastspeech2/>.

# TTS acoustic model

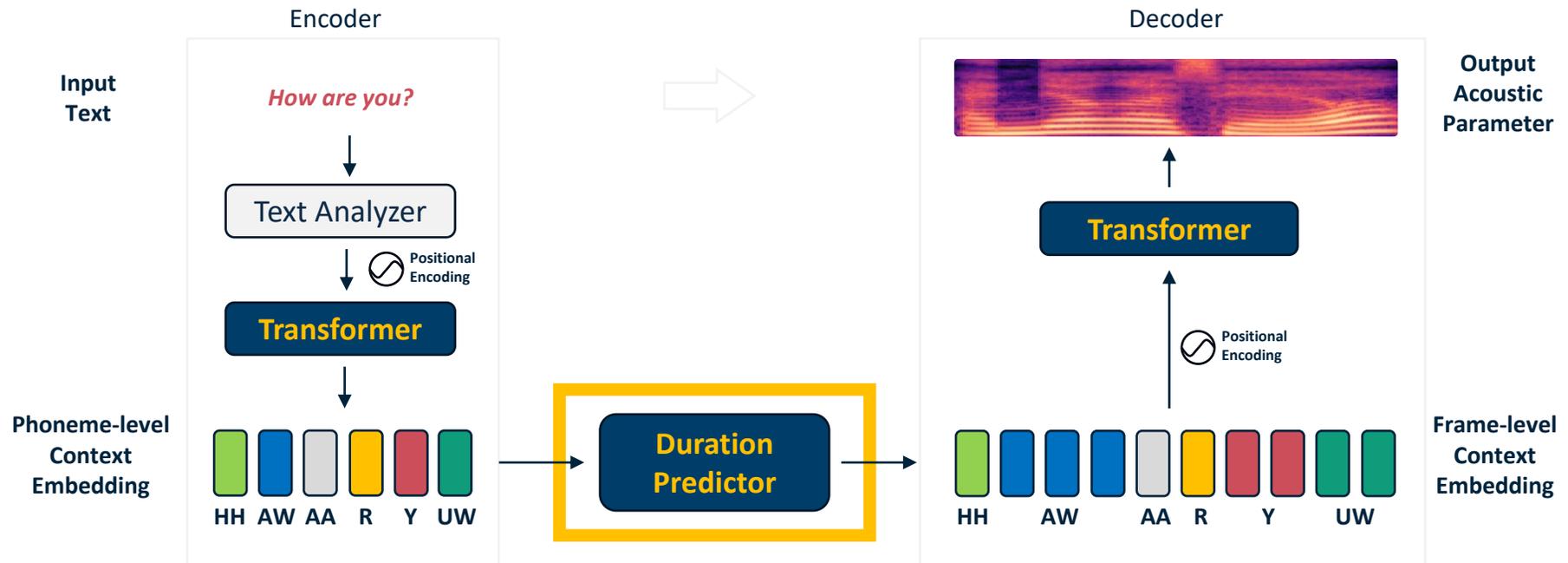
FastSpeech 2 (2020)



트랜스포머 기반의 인코더-디코더 사용

# TTS acoustic model

FastSpeech 2 (2020)



Duration predictor-based alignment

파라미터 복원을 병렬로(non-autoregressive) 처리함으로써 생성 속도를 개선

# Zero-shot Voice Cloning

---

Scaling TTS model

# Zero-shot voice cloning

## Recording constraint

	Conventional TTS	Voice cloning
Recording amount	> 30~60 min	< Few seconds
Speaking type	Script reading	Spontaneous speaking
Speaker	Professional voice actor	Non-professional
Recording amount	Clean studio	Anywhere
TTS quality	Natural	Unnatural

# Zero-shot voice cloning

## Recording constraint

	Conventional TTS	Voice cloning
Recording amount	> 30~60 min	< Few seconds
Speaking type	Script reading	Spontaneous speaking
Speaker	Professional voice actor	Non-professional
Recording amount	Clean studio	Anywhere
TTS quality	Natural	Unnatural

**Recording quality matters: Poor recording → TTS degradation**

# Zero-shot voice cloning

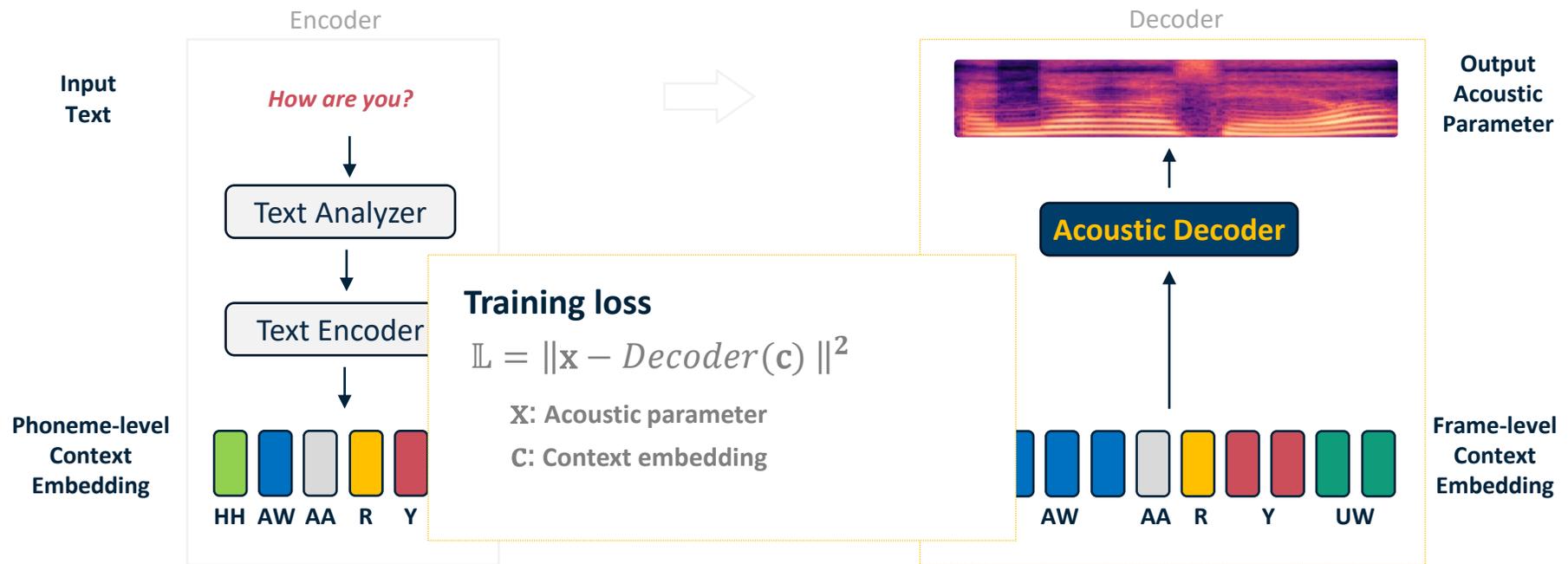
## Recording constraint

	Conventional TTS	Voice cloning
Recording amount	> 30~60 min	< Few seconds
Speaking type	Script reading	Spontaneous speaking
Speaker	Professional voice actor	Non-professional
Recording amount	Clean studio	Anywhere
TTS quality	Natural	<b>Very natural</b>

~~Recording quality matters: Poor recording → TTS degradation~~

# Zero-shot voice cloning

## Recall – Conventional TTS



The model directly learns characteristic of the target voice..

→ **Output quality** is heavily **dependent** on **target data**

# Zero-shot voice cloning

Voicebox (2023)

---

## 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\*

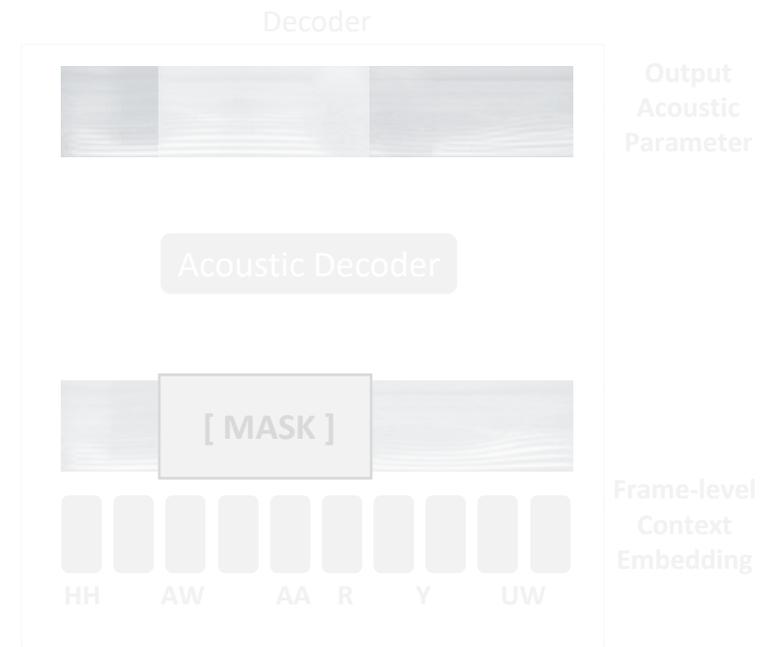
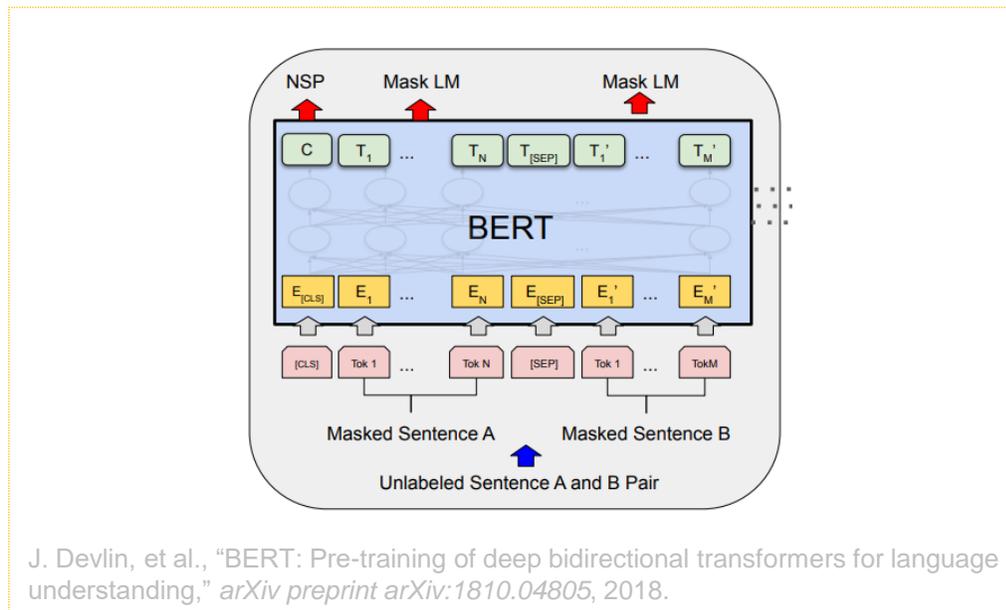
Fundamental AI Research (FAIR), Meta

### 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>.

# Zero-shot voice cloning

Key solution: Applying audio infilling task



Inspired by BERT's **masked language modeling**,

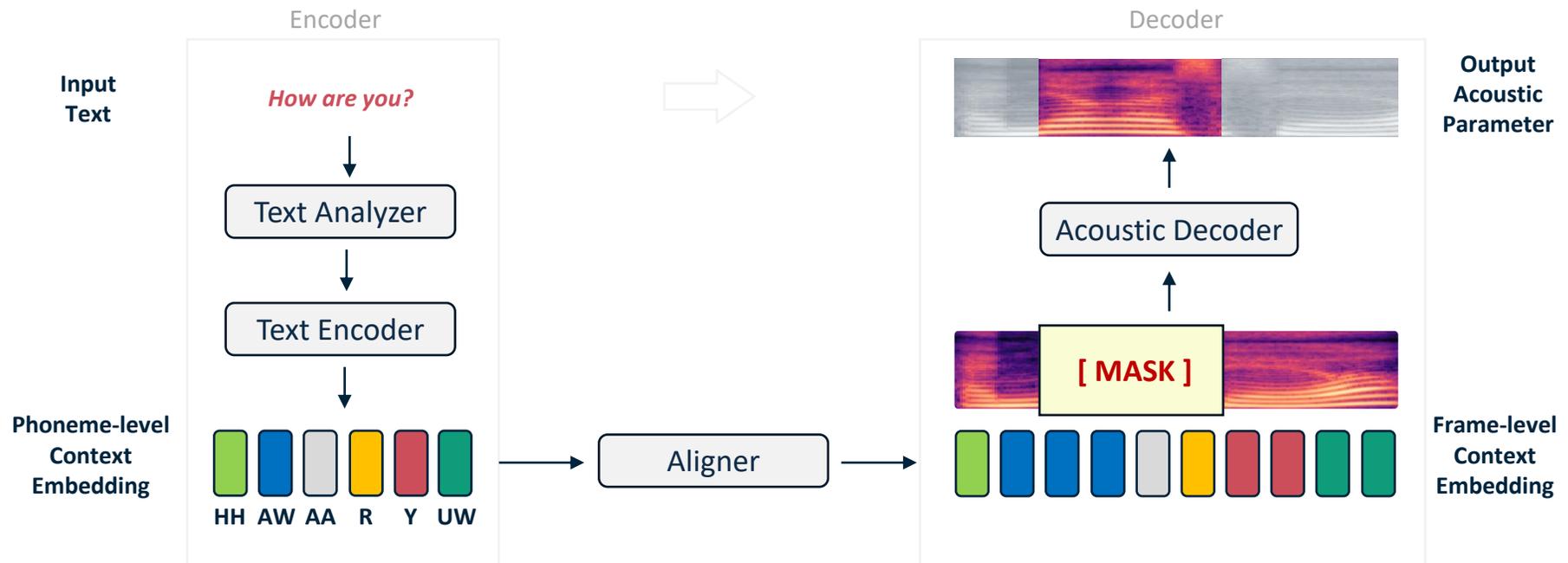
the model is trained to predict masked acoustic parameters using neighboring acoustic information



# Zero-shot voice cloning

Key solution: Applying audio infilling task

Training with **large-scale** speech corpora

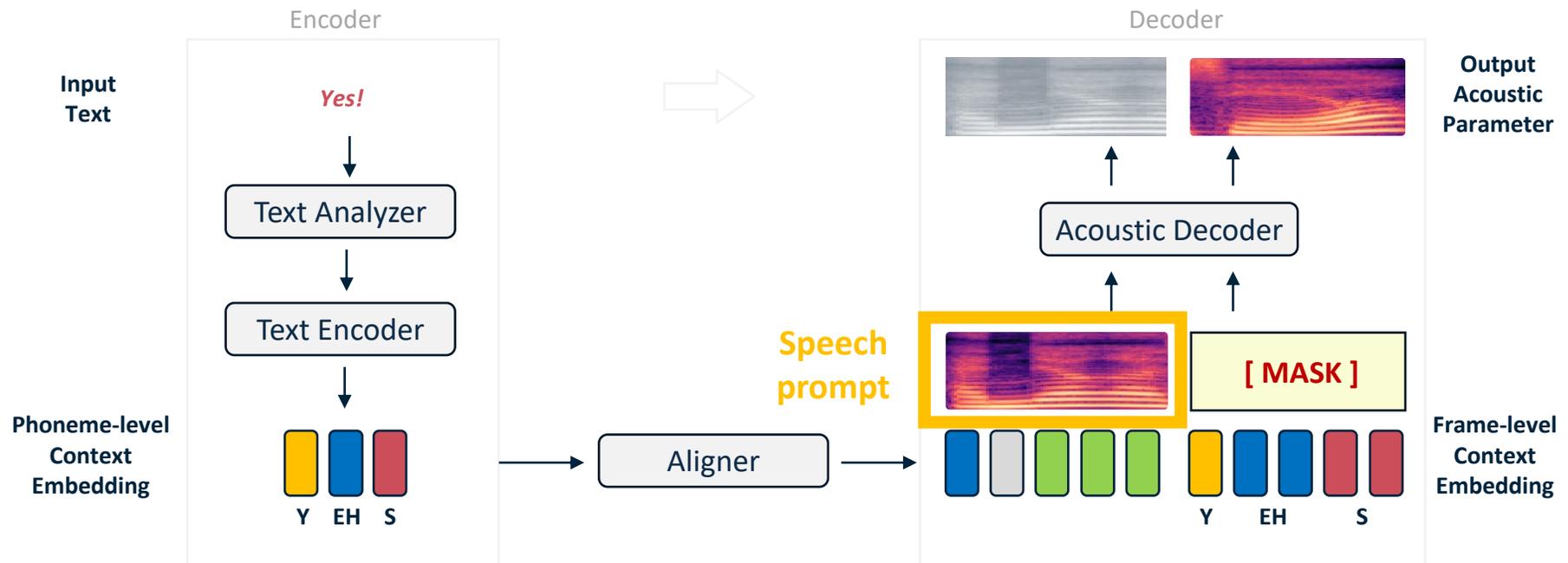


The model focuses on **relationship** between **adjacent acoustic parameters**,  
rather than reconstructing the target data

# Zero-shot voice cloning

Key solution: Applying audio infilling task

Inference with 5~10 seconds speech prompt



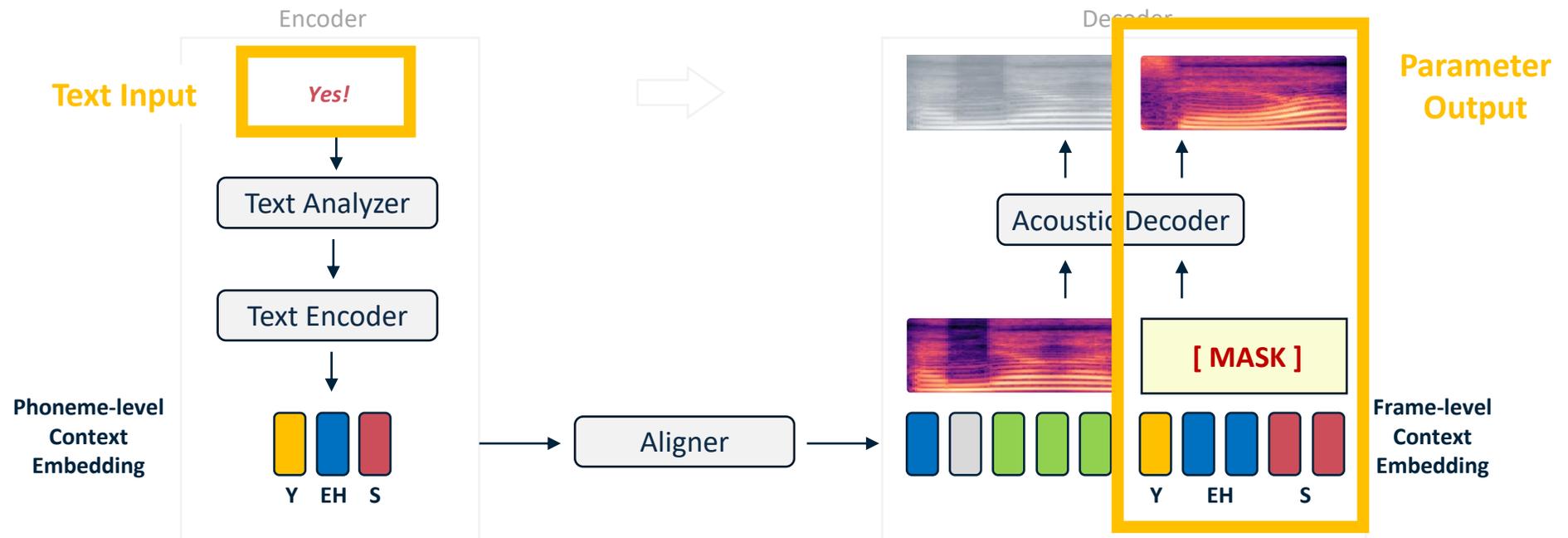
From the given **text** and **speech prompt**,

the model generates corresponding acoustic parameters

# Zero-shot voice cloning

Key solution: Applying audio infilling task

Inference with 5~10 seconds speech prompt

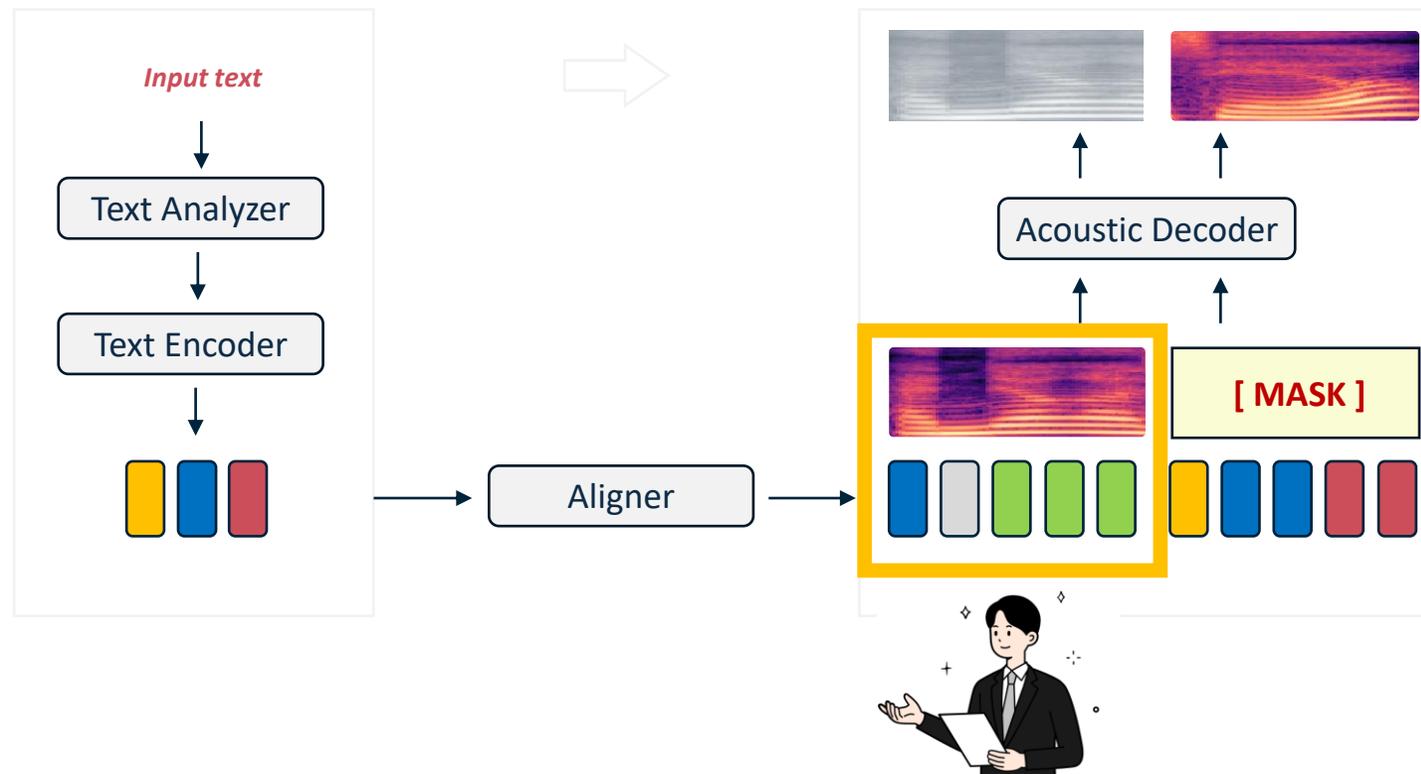


From the given **text** and **speech prompt**,  
the model **generates** corresponding **acoustic parameters**

# Zero-shot voice cloning

Key solution: Applying audio infilling task

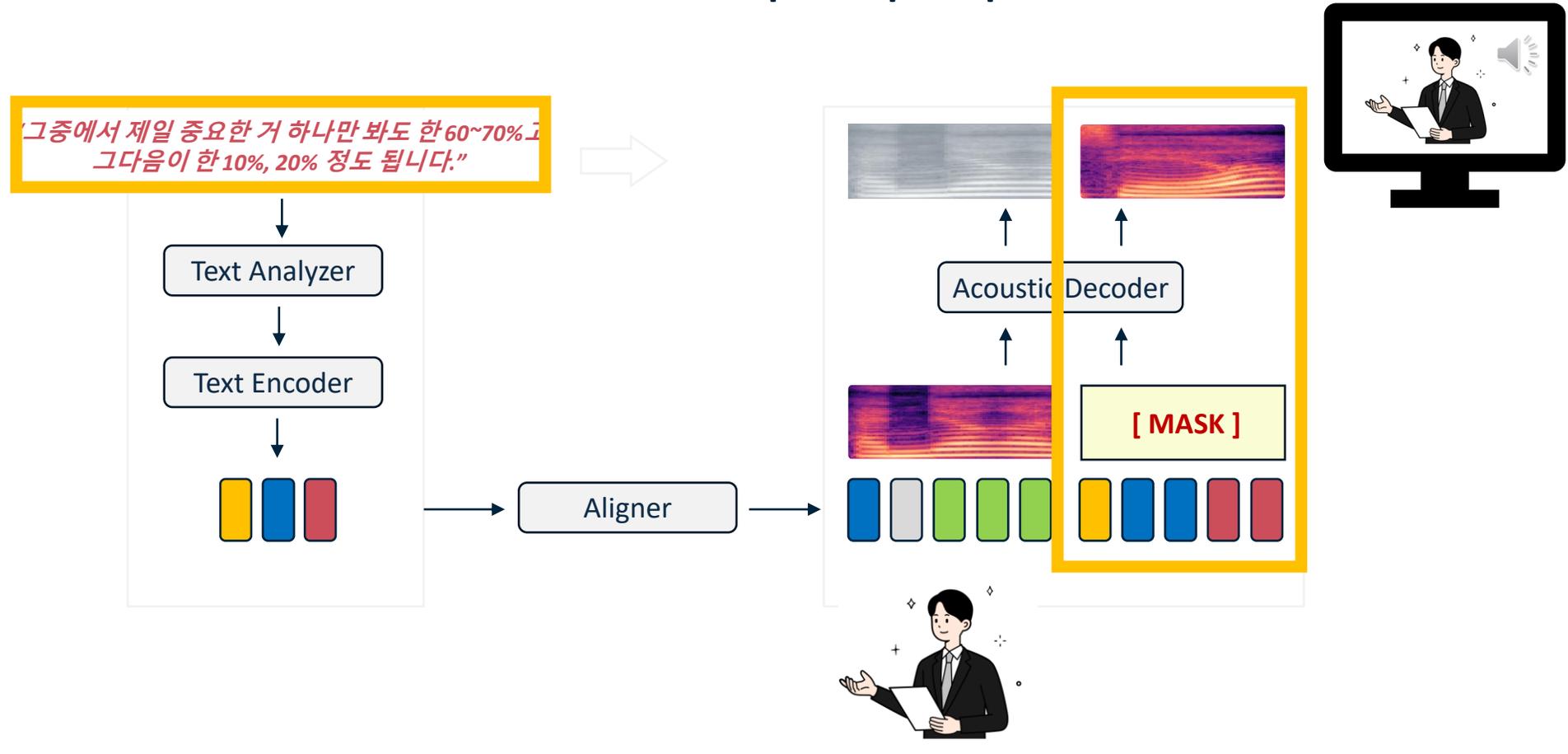
Inference with 5~10 seconds speech prompt



# Zero-shot voice cloning

Key solution: Applying audio infilling task

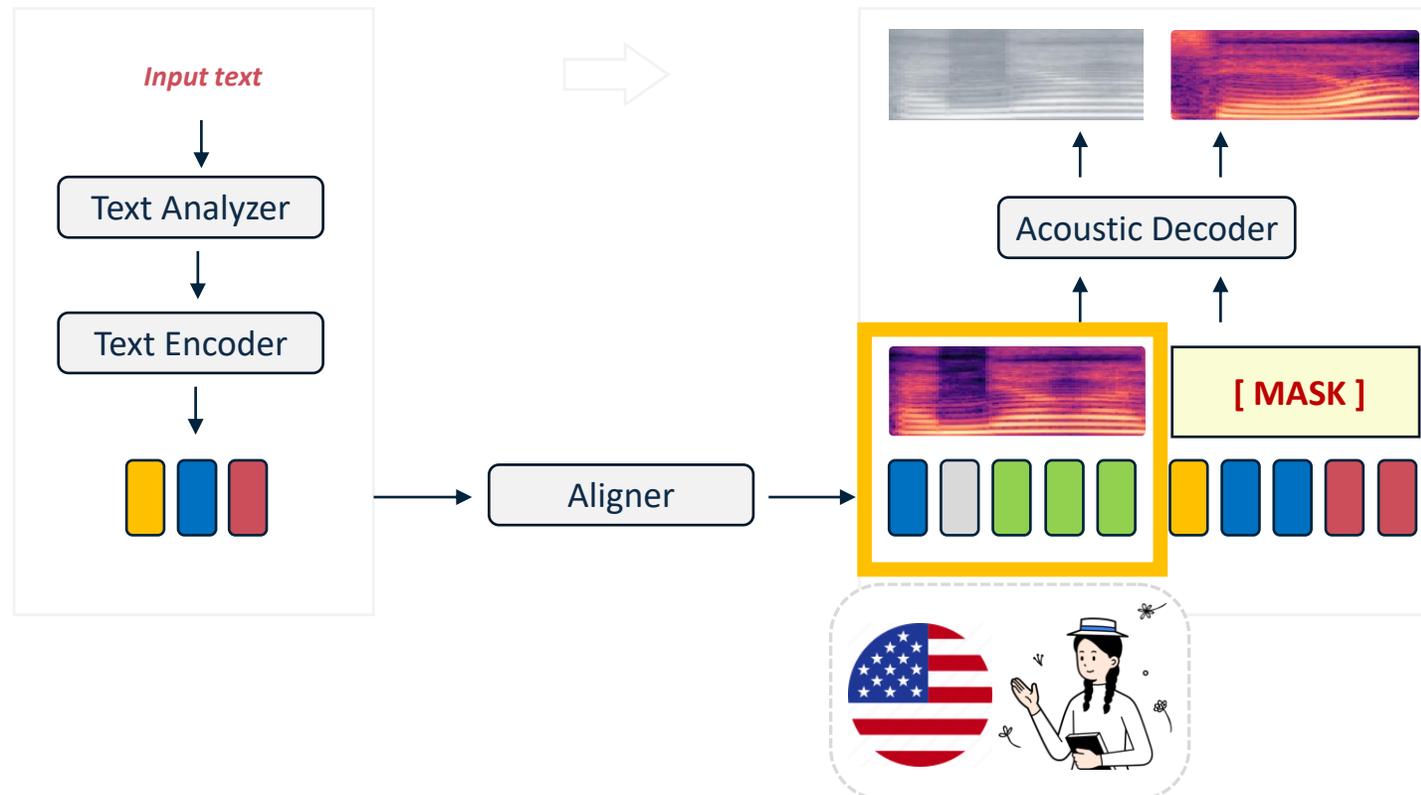
Inference with 5~10 seconds speech prompt



# Zero-shot voice cloning

Key solution: Applying audio infilling task

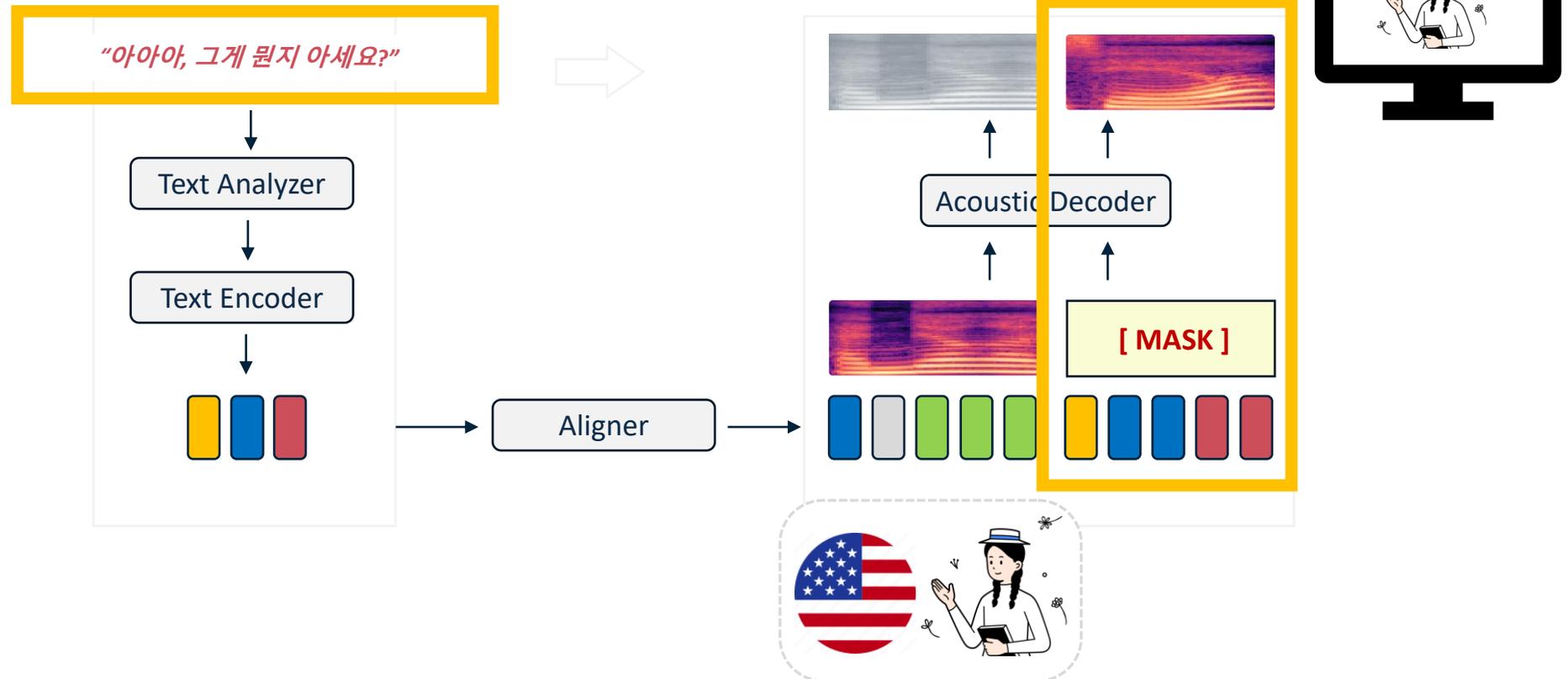
Inference with 5~10 seconds speech prompt



# Zero-shot voice cloning

Key solution: Applying audio infilling task

Inference with 5~10 seconds speech prompt



# Zero-shot voice cloning

## Overcoming the recording constraint

	Conventional TTS	Voice cloning
Speaker	Professional voice actor	Non-professional
Recording environment	Clean studio	Anywhere
Recording amount	> 30~60 min	< Few seconds
Speaking type	Clean studio	Anywhere
Model size	0.03B	0.41B
Inference speed	Real time x5 (CPU)	Real time x5 (GPU)
TTS quality	Natural	Very natural

# Zero-shot voice cloning

## Evaluations

	Conventional TTS	Voice cloning
Dataset	4 Korean speakers (2 male + 2 female)	
	~1h / speaker	4~8s / speaker
Speaker similarity (SECS)↑	68.0 %	78.3%
Intelligibility (CER)↓	1.8%	1.1%
Naturalness (MOS)↑	4.2	4.4

SECS; speaker embedding cosine similarity: 스피커 임베딩 벡터간의 유사도  
CER; character error rate: 입력 텍스트 ↔ 출력 음성의 ASR 결과(텍스트)간의 오류율  
MOS; mean opinion score: 전문가 청취평가(1~5 scale)

# Zero-shot Voice Cloning

---

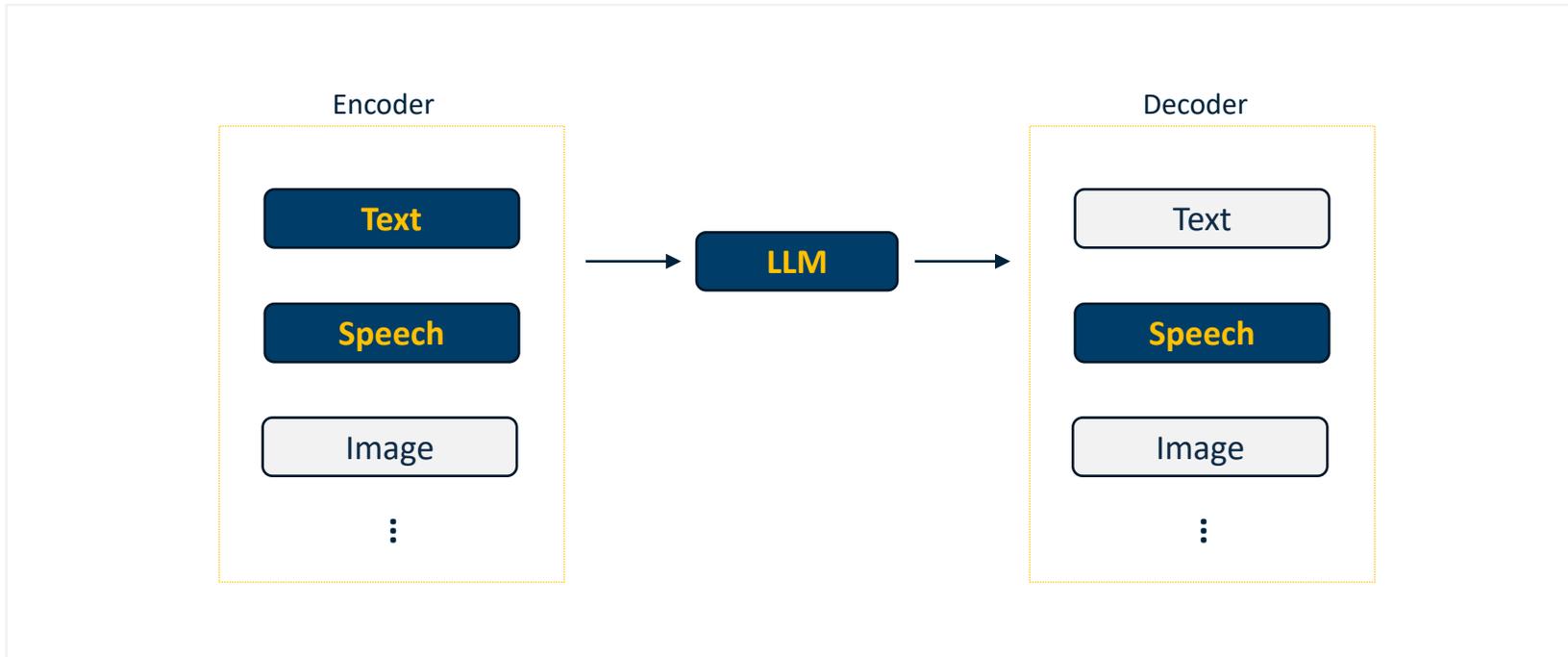
With **large language models**



Source of the original video: <https://www.youtube.com/watch?v=c2DFg53Zhvw>  
"Live demo of GPT-4o realtime translation"

# Zero-shot voice cloning

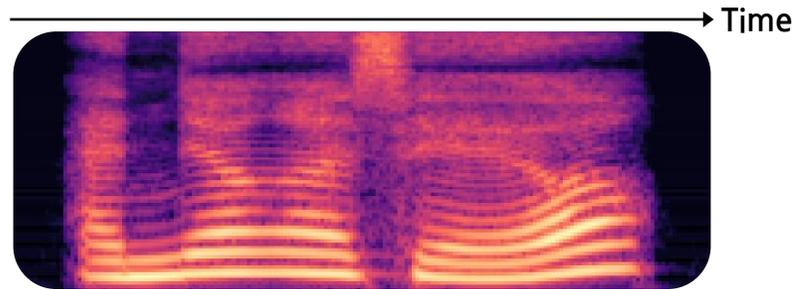
Toward Voice Agent



LLM 능력을 음성 도메인까지 확장할 수 있을까?

# Zero-shot voice cloning

Recall - melspectrogram



Acoustic parameters..?

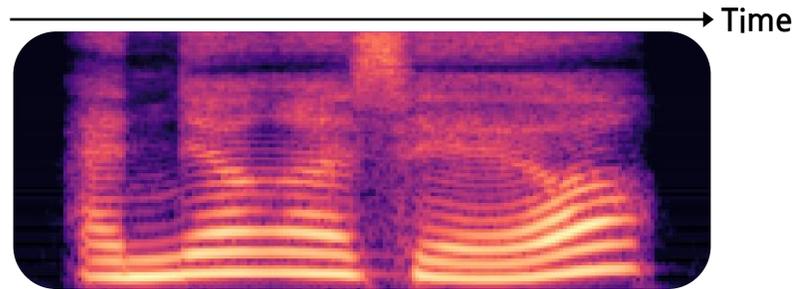
Speaker-specific attributes  
(tone, volume, timbre, speaking rate, ...)

**Continuous representation**

→ Discrete token 을 사용하는 LLM 모델링에 적합하지 않음

# Zero-shot voice cloning

Recall - melspectrogram



Acoustic parameters..?

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

**Continuous representation**

→ **Discrete token** 을 사용하는 LLM 모델링에 적합하지 않음

# Zero-shot voice cloning

## CosyVoice 2 (2024)

---

### CosyVoice 2: Scalable Streaming Speech Synthesis with Large Language Models

---

Zhihao Du, Yuxuan Wang, Qian Chen, Xian Shi, Xiang Lv, Tianyu Zhao, Zhifu Gao  
Yexin Yang, Changfeng Gao, Hui Wang, Fan Yu, Huadai Liu, Zhengyan Sheng  
Yue Gu, Chong Deng, Wen Wang, Shiliang Zhang, Zhijie Yan, Jingren Zhou\*

Alibaba Group, China

{neo.dzh,sly.zsl}@alibaba-inc.com

#### Abstract

In our previous work, we introduced CosyVoice, a multilingual speech synthesis model based on supervised discrete speech tokens. By employing progressive semantic decoding with two popular generative models, language models (LMs) and Flow Matching, CosyVoice demonstrated high prosody naturalness, content consistency, and speaker similarity in speech in-context learning. Recently, significant progress has been made in multi-modal large language models (LLMs), where the response latency and real-time factor of speech synthesis play a crucial role in the interactive experience. Therefore, in this report, we present an improved streaming speech synthesis model, CosyVoice 2, which incorporates comprehensive and systematic optimizations. Specifically, we introduce finite-scalar quantization to improve the codebook utilization of speech tokens. For the text-speech LM, we streamline the model architecture to allow direct use of a pre-trained LLM as the backbone. In addition, we develop a chunk-aware causal flow matching model to support various synthesis scenarios, enabling both streaming and non-streaming synthesis within a single model. By training on a large-scale multilingual dataset, CosyVoice 2 achieves human-parity naturalness, minimal response latency, and virtually lossless synthesis quality in the streaming mode. We invite readers to listen to the demos at <https://funaudiollm.github.io/cosyvoice2>.

# Zero-shot voice cloning

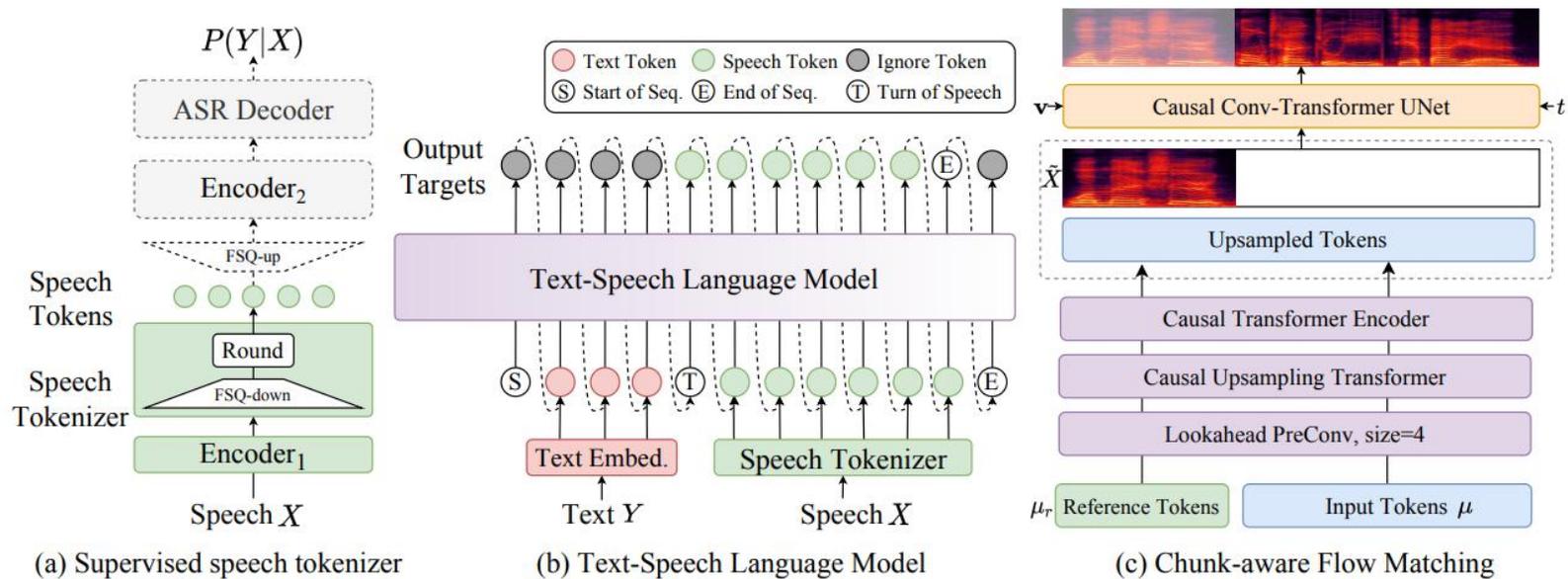
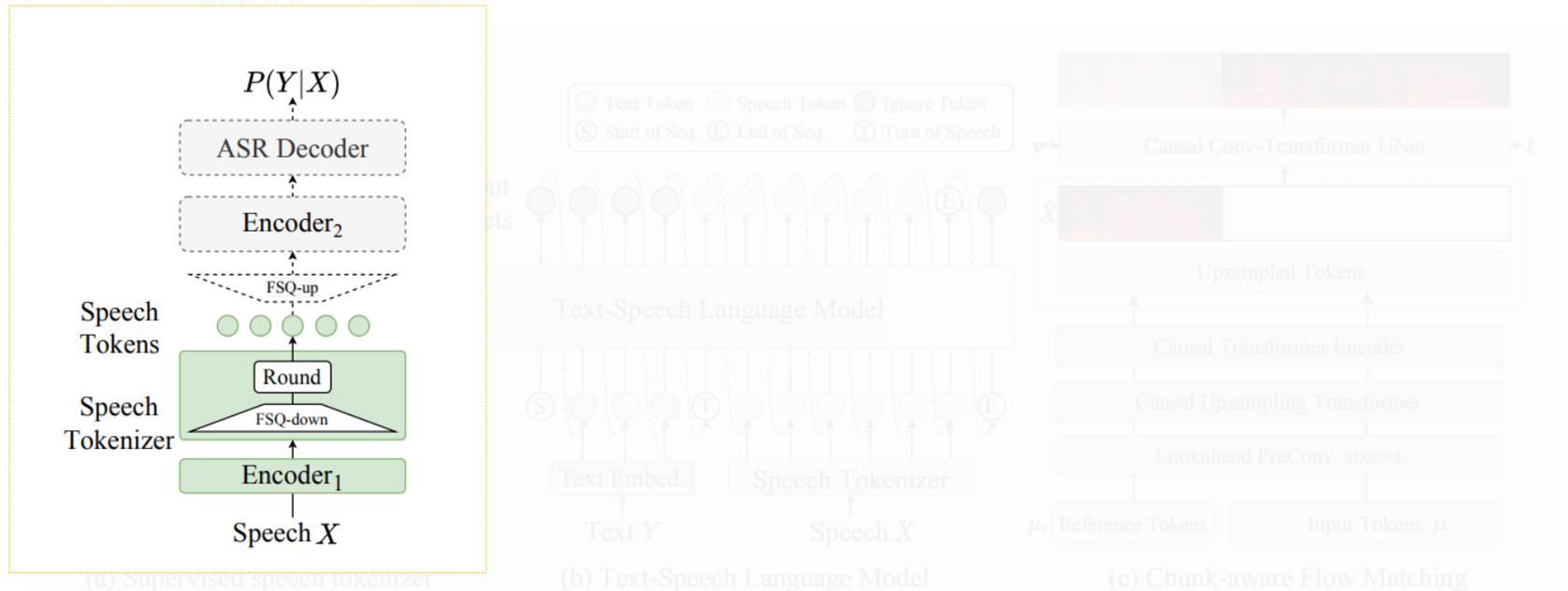


Figure 1: An overview of CosyVoice 2. (a) demonstrates the supervised speech tokenizer, where dashed modules are only used at the training stage. (b) is a unified text-speech language model for streaming and non-streaming synthesis. Dashed lines indicate the autoregressive decoding at the inference stage. (c) illustrates the causal flow matching model conditioning on a speaker embedding  $\mathbf{v}$ , semantic tokens  $\mu$ , masked speech features  $\tilde{X}$  and intermediate state  $X_t$  at timestep  $t$  on the probabilistic density path.

# Zero-shot voice cloning

## Speech tokenizer

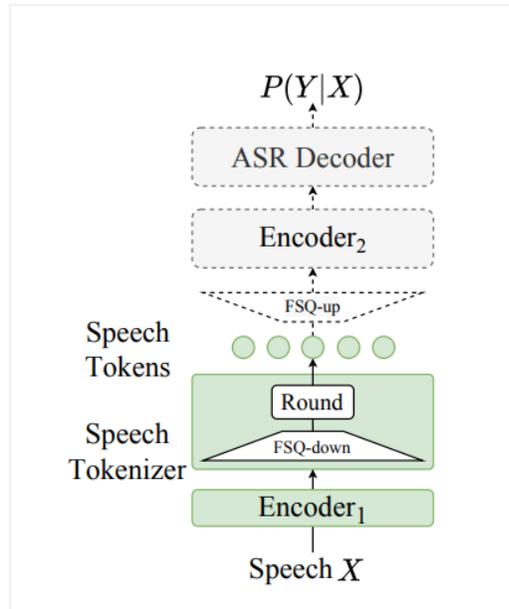


### Speech $\rightarrow$ Speech Token

of CosyVoice 2. (a) demonstrates the supervised speech tokenizer, where dashed modules are only used at the training stage. (b) is a unified text-speech language model for streaming and non-streaming synthesis. Dashed lines indicate the autoregressive decoding at the inference stage. (c) illustrates the causal flow matching model conditioning on a speaker embedding  $v$ , semantic tokens  $\mu$ , masked speech features  $\tilde{X}$  and intermediate state  $X_t$  at timestep  $t$  on the probabilistic density path.

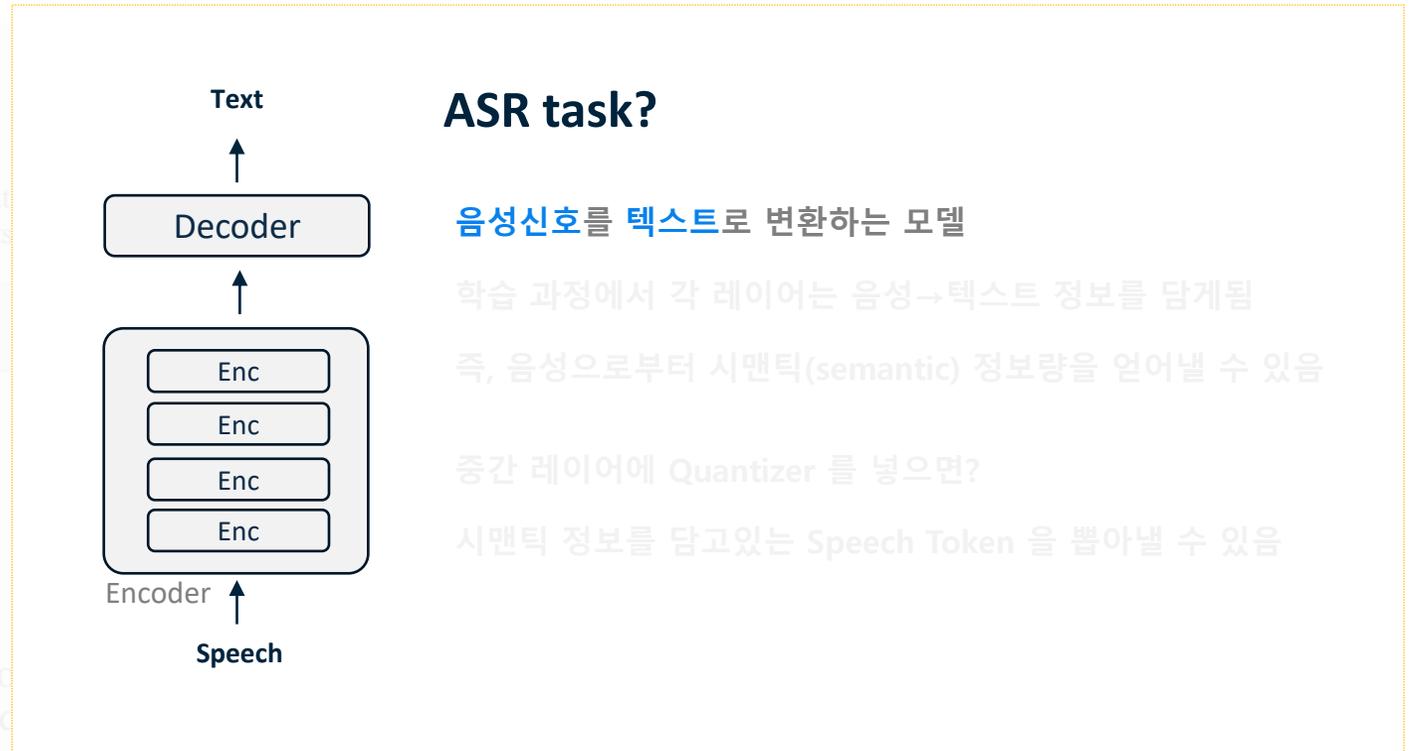
# Zero-shot voice cloning

## Speech tokenizer



(a) Supervised speech tokenizer

Figure 1: An overview of Co...  
dashed modules are only used...  
streaming and non-streaming synthesis. Dashed lines indicate the autoregressive decoding at the inference stage. (c) illustrates the causal flow matching model conditioning on a speaker embedding  $\mathbf{v}$ , semantic tokens  $\mu$ , masked speech features  $\tilde{X}$  and intermediate state  $X_t$  at timestep  $t$  on the probabilistic density path.



### ASR task?

음성신호를 텍스트로 변환하는 모델

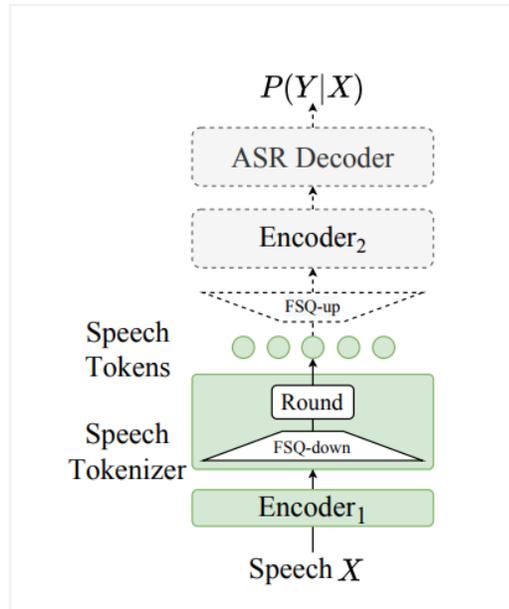
학습 과정에서 각 레이어는 음성→텍스트 정보를 담게됨  
즉, 음성으로부터 시맨틱(semantic) 정보량을 얻어낼 수 있음

중간 레이어에 Quantizer 를 넣으면?

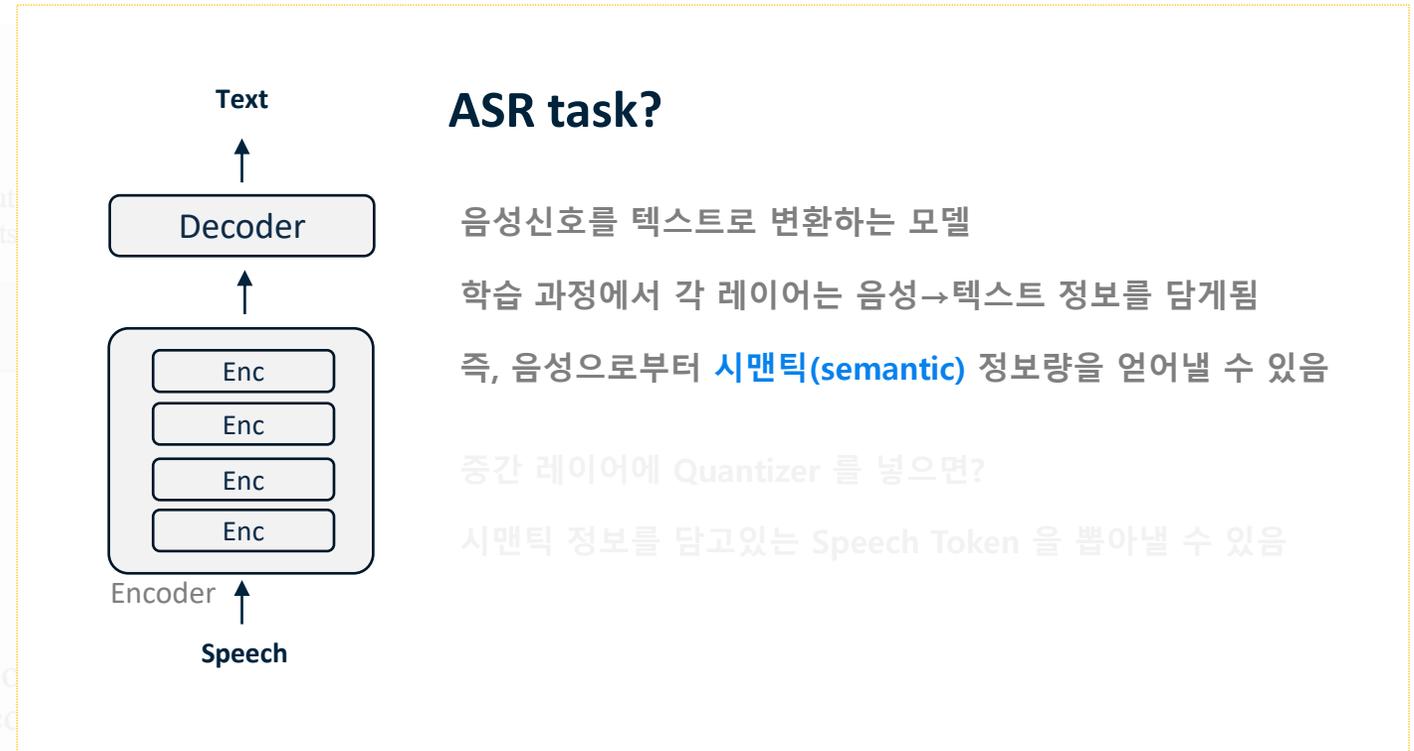
시맨틱 정보를 담고있는 Speech Token 을 뽑아낼 수 있음

# Zero-shot voice cloning

## Speech tokenizer



(a) Supervised speech tokenizer



### ASR task?

음성신호를 텍스트로 변환하는 모델

학습 과정에서 각 레이어는 음성→텍스트 정보를 담게 됨

즉, 음성으로부터 **시맨틱(semantic)** 정보량을 얻어낼 수 있음

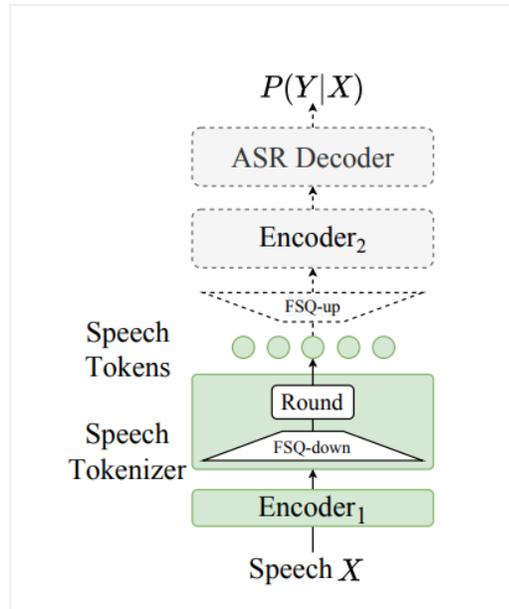
중간 레이어에 Quantizer 를 넣으면?

시맨틱 정보를 담고있는 Speech Token 을 뽑아낼 수 있음

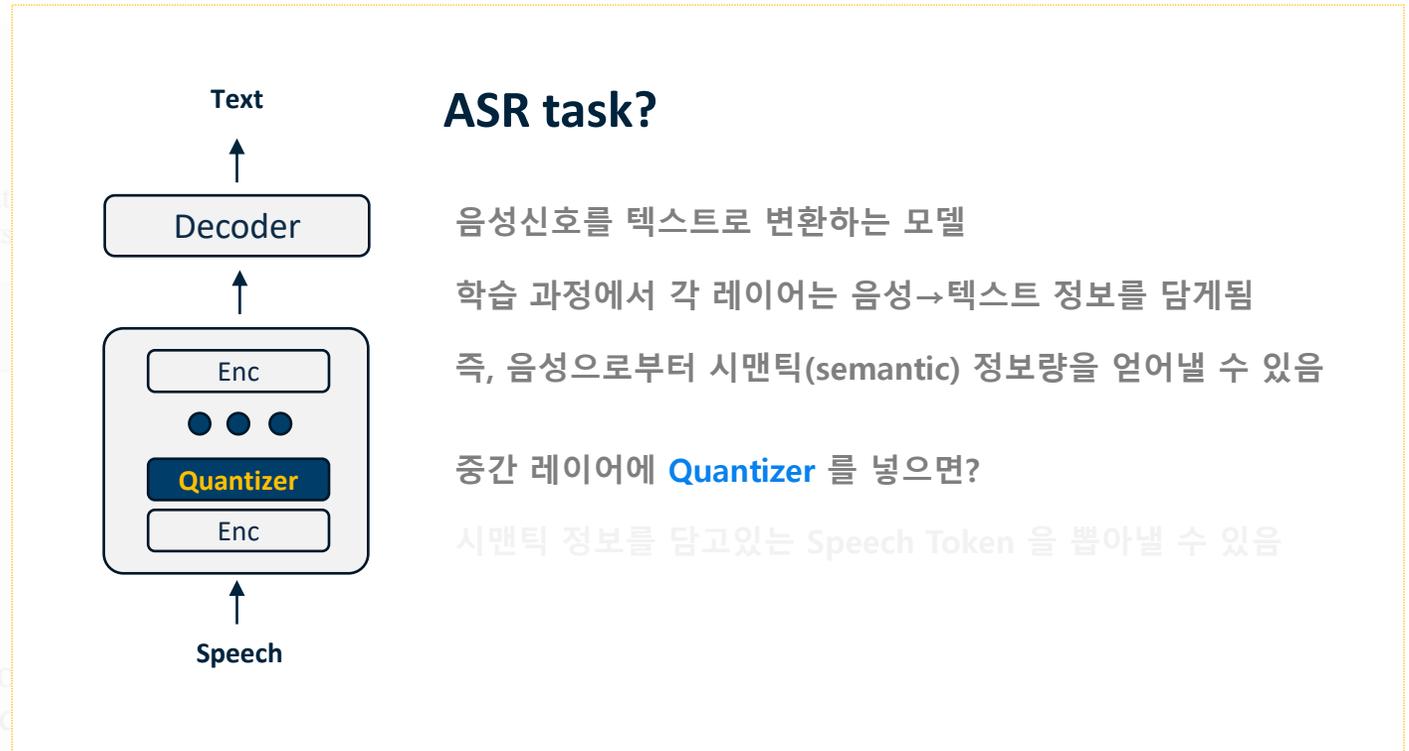
Figure 1: An overview of CoCo. Dashed modules are only used for streaming and non-streaming synthesis. Dashed lines indicate the autoregressive decoding at the inference stage. (c) illustrates the causal flow matching model conditioning on a speaker embedding  $\mathbf{v}$ , semantic tokens  $\mu$ , masked speech features  $\tilde{X}$  and intermediate state  $X_t$  at timestep  $t$  on the probabilistic density path.

# Zero-shot voice cloning

## Speech tokenizer



(a) Supervised speech tokenizer



### ASR task?

음성신호를 텍스트로 변환하는 모델

학습 과정에서 각 레이어는 음성→텍스트 정보를 담게 됨

즉, 음성으로부터 시맨틱(semantic) 정보량을 얻어낼 수 있음

중간 레이어에 **Quantizer** 를 넣으면?

시맨틱 정보를 담고있는 **Speech Token** 을 뽑아낼 수 있음

Figure 1: An overview of CoCo. (a) shows the supervised speech tokenizer. (b) shows the ASR task model. Dashed modules are only used for streaming and non-streaming synthesis. Dashed lines indicate the autoregressive decoding at the inference stage. (c) illustrates the causal flow matching model conditioning on a speaker embedding  $\mathbf{v}$ , semantic tokens  $\mu$ , masked speech features  $\tilde{X}$  and intermediate state  $X_t$  at timestep  $t$  on the probabilistic density path.

# Zero-shot voice cloning

## Speech tokenizer

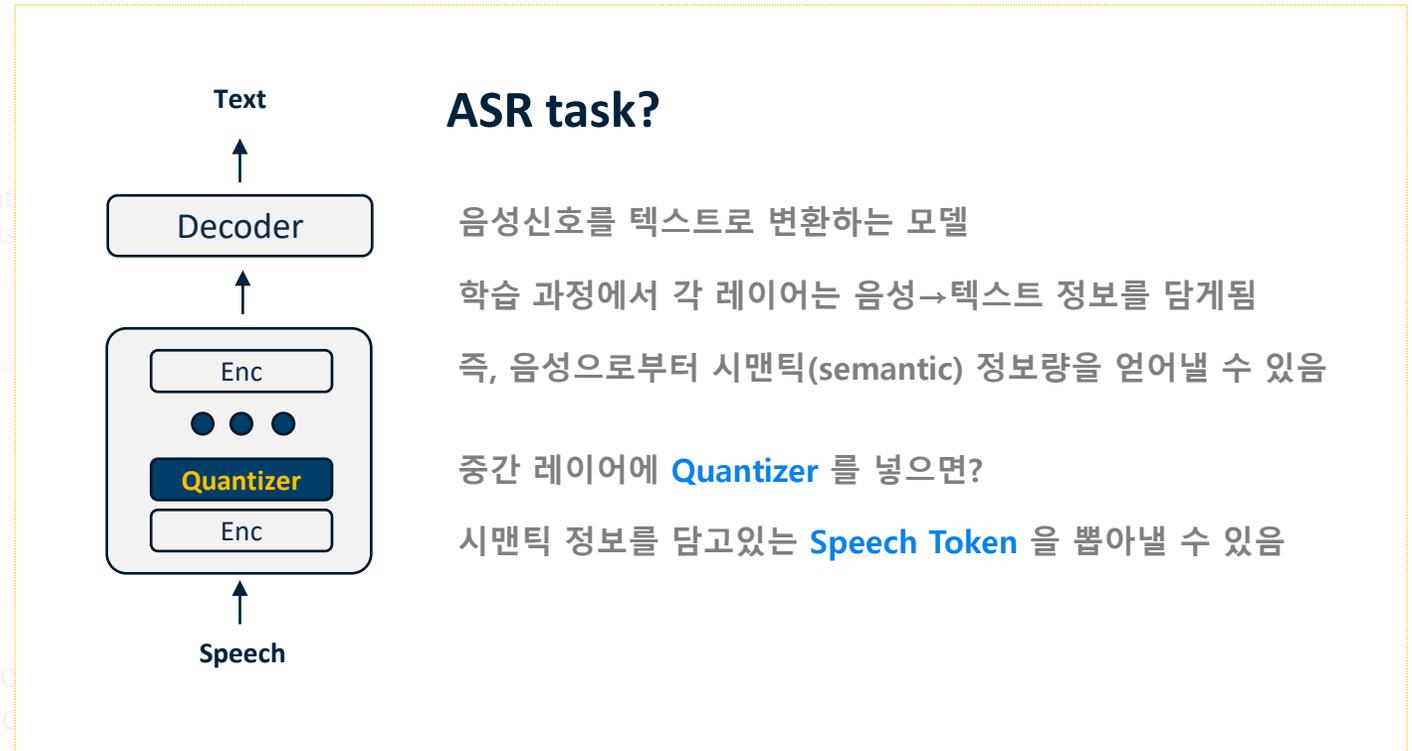
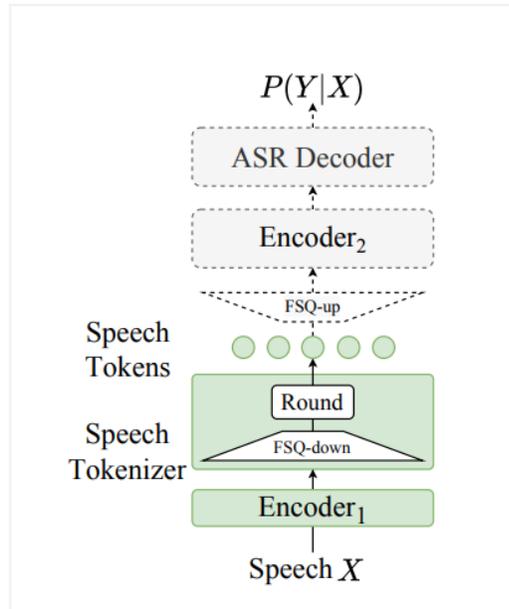
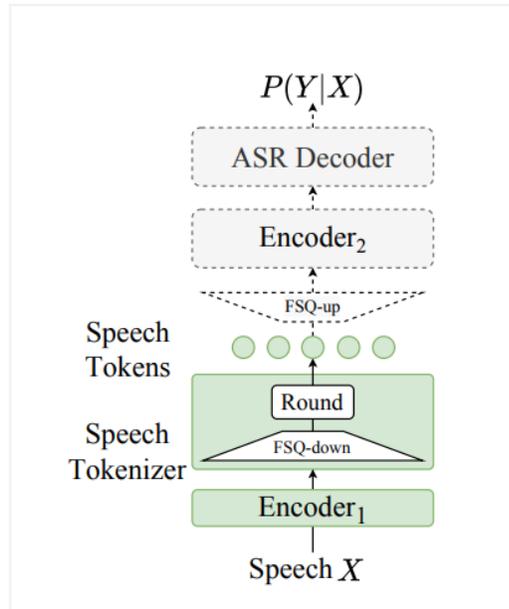


Figure 1: An overview of Co...  
dashed modules are only used...  
streaming and non-streaming synthesis. Dashed lines indicate the autoregressive decoding at the inference stage. (c) illustrates the causal flow matching model conditioning on a speaker embedding  $\mathbf{v}$ , semantic tokens  $\mu$ , masked speech features  $\tilde{X}$  and intermediate state  $X_t$  at timestep  $t$  on the probabilistic density path.

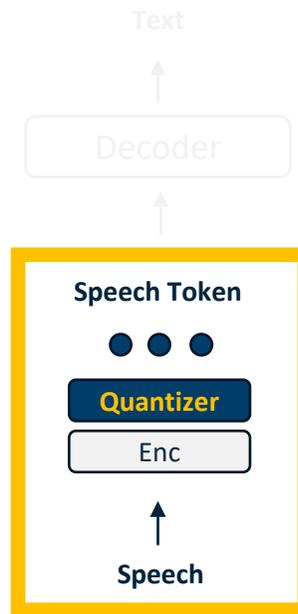
# Zero-shot voice cloning

## Speech tokenizer



(a) Supervised speech tokenizer

Figure 1: An overview of CoCo. Dashed modules are only used for streaming and non-streaming synthesis. Dashed lines indicate the autoregressive decoding at the inference stage. (c) illustrates the causal flow matching model conditioning on a speaker embedding  $\mathbf{v}$ , semantic tokens  $\mu$ , masked speech features  $\tilde{X}$  and intermediate state  $X_t$  at timestep  $t$  on the probabilistic density path.



### ASR task?

음성신호를 텍스트로 변환하는 모델

학습 과정에서 각 레이어는 음성→텍스트 정보를 담게 됨

즉, 음성으로부터 시맨틱(semantic) 정보량을 얻어낼 수 있음

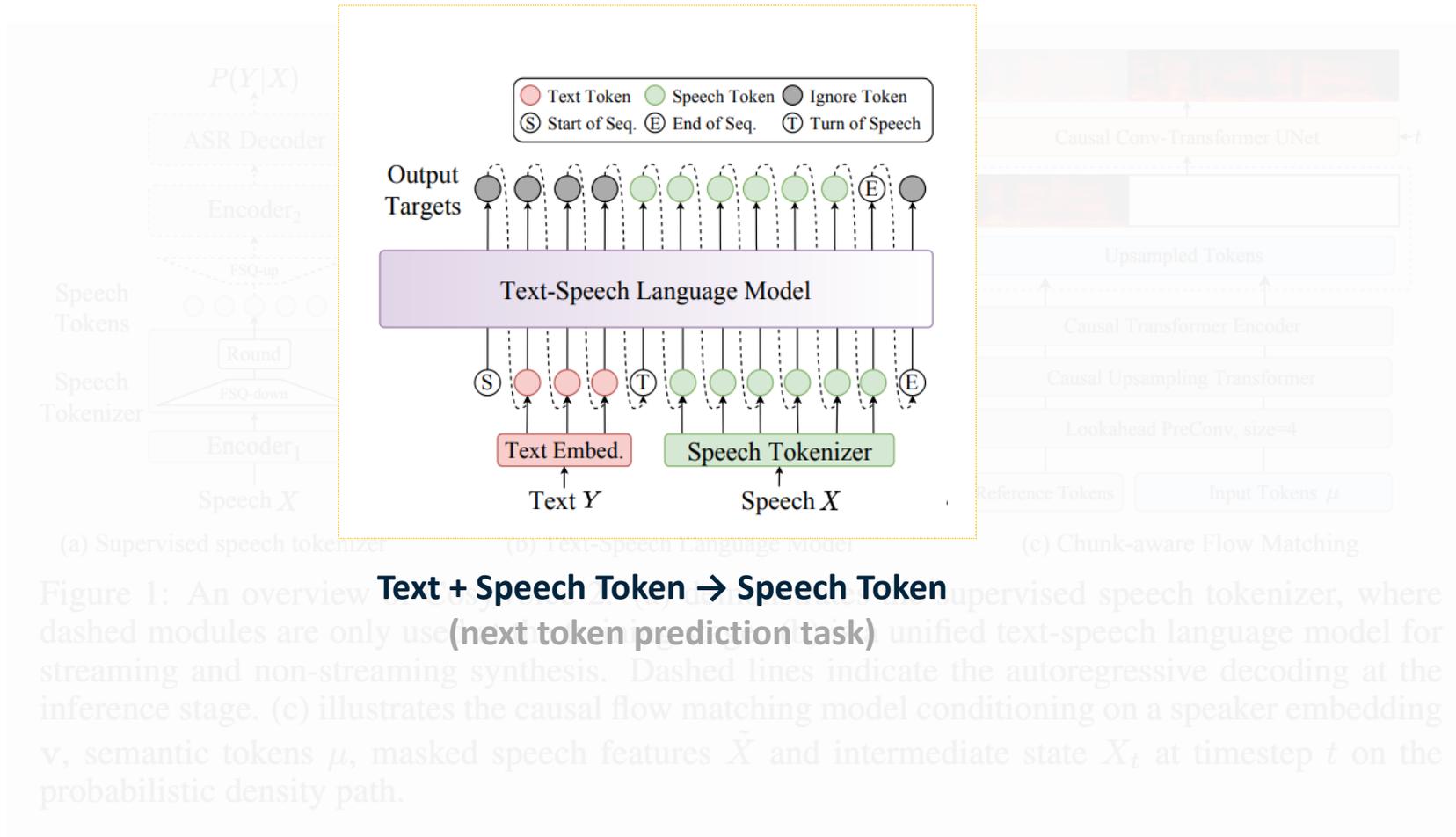
중간 레이어에 **Quantizer** 를 넣으면?

시맨틱 정보를 담고있는 **Speech Token** 을 뽑아낼 수 있음

### Speech Tokenizer

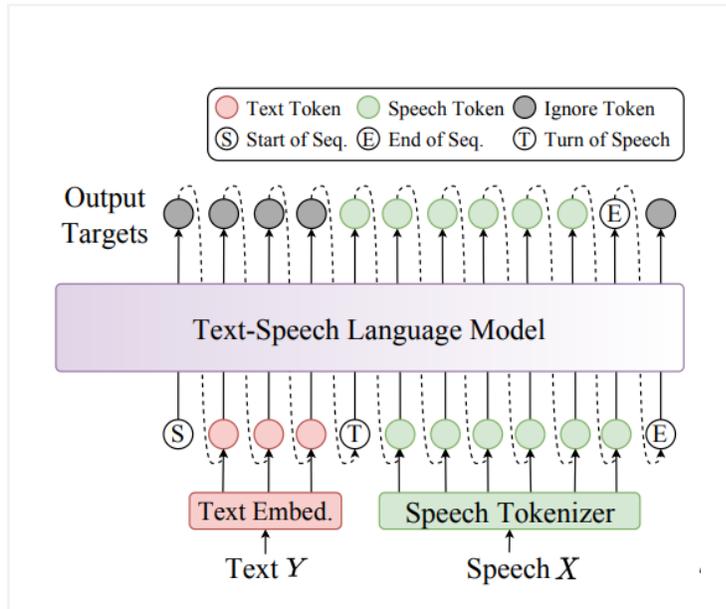
# Zero-shot voice cloning

## Autoregressive language model



# Zero-shot voice cloning

## Autoregressive language model



(a) Supervised speech tokenizer

Figure 1: An overview of Co...  
dashed modules are only used  
streaming and non-streaming  
inference stage. (c) illustrates  
 $v$ , semantic tokens  $\mu$ , masked  
probabilistic density path.

### Next token prediction task?

입력 토큰 다음에 올 토큰의 확률 분포를 예측함으로써 학습 데이터의 통계적 패턴을 학습

### TTS design - training

input = "<|sos|>{text token}<|tos|>{speech token} <|eos|>"

텍스트 토큰 뒤에 음성 토큰을 연결해 입력 시퀀스를 구성하고,

모델이 텍스트 의미를 해석한 뒤 해당 음성 토큰을 순차적으로 재구성하도록 학습

### TTS design - inference (in-context learning)

input = "<|sos|>{prompt text token}{input text token}<|tos|>{prompt speech token}"

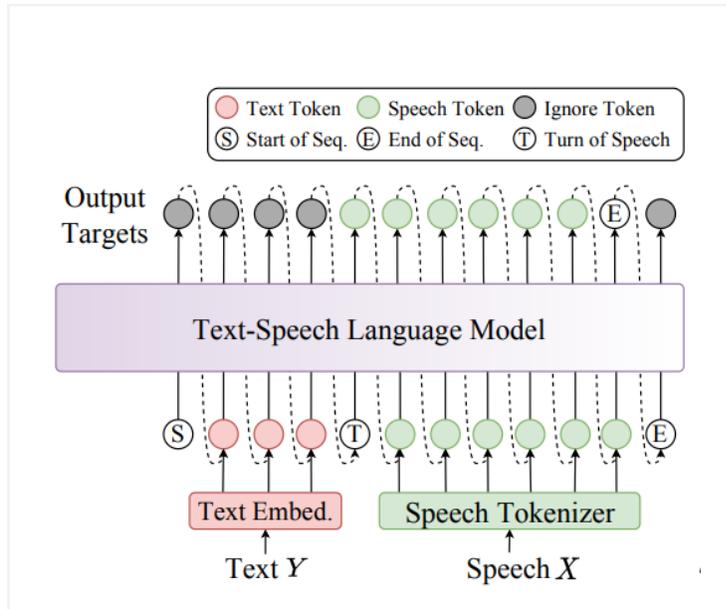
Output = input + "{output speech token}<|eos|>"

프롬프트(텍스트, 음성) 토큰과 입력 텍스트 토큰의 의미를 해석한 뒤

출력 음성 토큰을 순차적으로 추론

# Zero-shot voice cloning

## Autoregressive language model



(a) Supervised speech tokenizer

Figure 1: An overview of Co...  
dashed modules are only used  
streaming and non-streaming  
inference stage. (c) illustrates  
 $v$ , semantic tokens  $\mu$ , masked  
probabilistic density path.

### Next token prediction task?

입력 토큰 다음에 올 토큰의 확률 분포를 예측함으로써 학습 데이터의 통계적 패턴을 학습

### TTS design - training

```
input = "<|sos|>{text token}<|tos|>{speech token}<|eos|>"
```

텍스트 토큰 뒤에 음성 토큰을 연결해 입력 시퀀스를 구성하고,

모델이 텍스트 의미를 해석한 뒤 해당 음성 토큰을 순차적으로 재구성하도록 학습

### TTS design - inference (in-context learning)

```
input = "<|sos|>{prompt text token}{input text token}<|tos|>{prompt speech token}"
```

```
Output = input + "{output speech token}<|eos|>"
```

프롬프트(텍스트, 음성) 토큰과 입력 텍스트 토큰의 의미를 해석한 뒤

출력 음성 토큰을 순차적으로 추론

# Zero-shot voice cloning

## Autoregressive language model

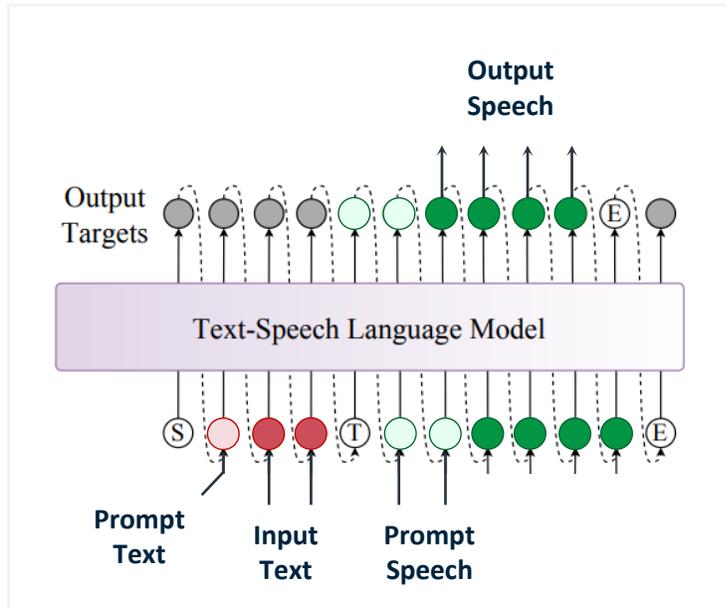


Figure 1: An overview of Co...  
dashed modules are only used  
streaming and non-streaming  
inference stage. (c) illustrates  
 $v$ , semantic tokens  $\mu$ , masked  
probabilistic density path.

### Next token prediction task?

입력 토큰 다음에 올 토큰의 확률 분포를 예측함으로써 학습 데이터의 통계적 패턴을 학습

### TTS design - training

```
input = "<|sos|>{text token}<|tos|>{speech token}<|eos|>"
```

텍스트 토큰 뒤에 음성 토큰을 연결해 입력 시퀀스를 구성하고,

모델이 텍스트 의미를 해석한 뒤 해당 음성 토큰을 순차적으로 재구성하도록 학습

### TTS design - inference (in-context learning)

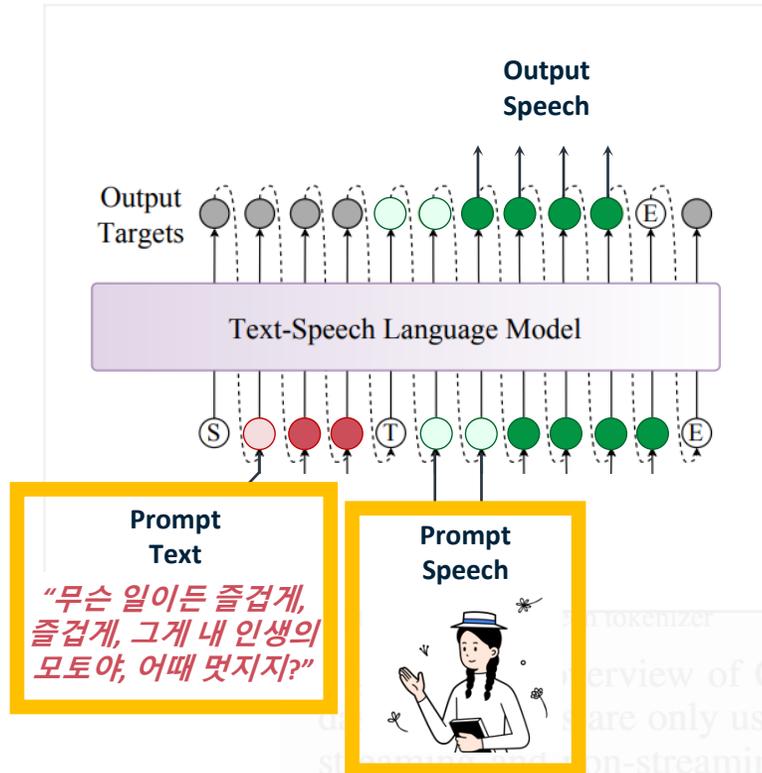
```
input = "<|sos|>{prompt text token}{input text token}<|tos|>{prompt speech token}"  
Output = input + "{output speech token}<|eos|>"
```

프롬프트(텍스트, 음성) 토큰과 입력 텍스트 토큰의 의미를 해석한 뒤

출력 음성 토큰을 순차적으로 추론

# Zero-shot voice cloning

## Autoregressive language model



### Next token prediction task?

입력 토큰 다음에 올 토큰의 확률 분포를 예측함으로써 학습 데이터의 통계적 패턴을 학습

### TTS design - training

```
input = "<|sos|>{text token}<|tos|>{speech token}<|eos|>"
```

텍스트 토큰 뒤에 음성 토큰을 연결해 입력 시퀀스를 구성하고,

모델이 텍스트 의미를 해석한 뒤 해당 음성 토큰을 순차적으로 재구성하도록 학습

### TTS design - inference (in-context learning)

```
input = "<|sos|>{prompt text token}{input text token}<|tos|>{prompt speech token}"
```

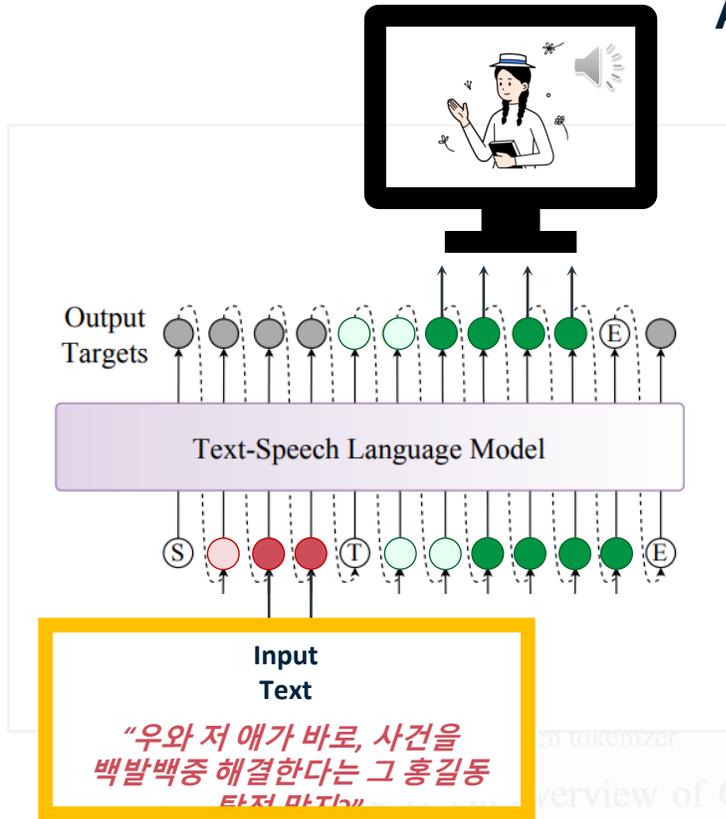
```
Output = input + "{output speech token}<|eos|>"
```

프롬프트(텍스트, 음성) 토큰과 입력 텍스트 토큰의 의미를 해석한 뒤

출력 음성 토큰을 순차적으로 추론

# Zero-shot voice cloning

## Autoregressive language model



### Next token prediction task?

입력 토큰 다음에 올 토큰의 확률 분포를 예측함으로써 학습 데이터의 통계적 패턴을 학습

### TTS design - training

```
input = "<|sos|>{text token}<|tos|>{speech token}<|eos|>"
```

텍스트 토큰 뒤에 음성 토큰을 연결해 입력 시퀀스를 구성하고,

모델이 텍스트 의미를 해석한 뒤 해당 음성 토큰을 순차적으로 재구성하도록 학습

### TTS design - inference (in-context learning)

```
input = "<|sos|>{prompt text token}{input text token}<|tos|>{prompt speech token}"
```

```
Output = input + "{output speech token}<|eos|>"
```

프롬프트(텍스트, 음성) 토큰과 **입력 텍스트 토큰**의 의미를 해석한 뒤

**출력 음성 토큰**을 순차적으로 추론

# Zero-shot voice cloning

## Speech decoder

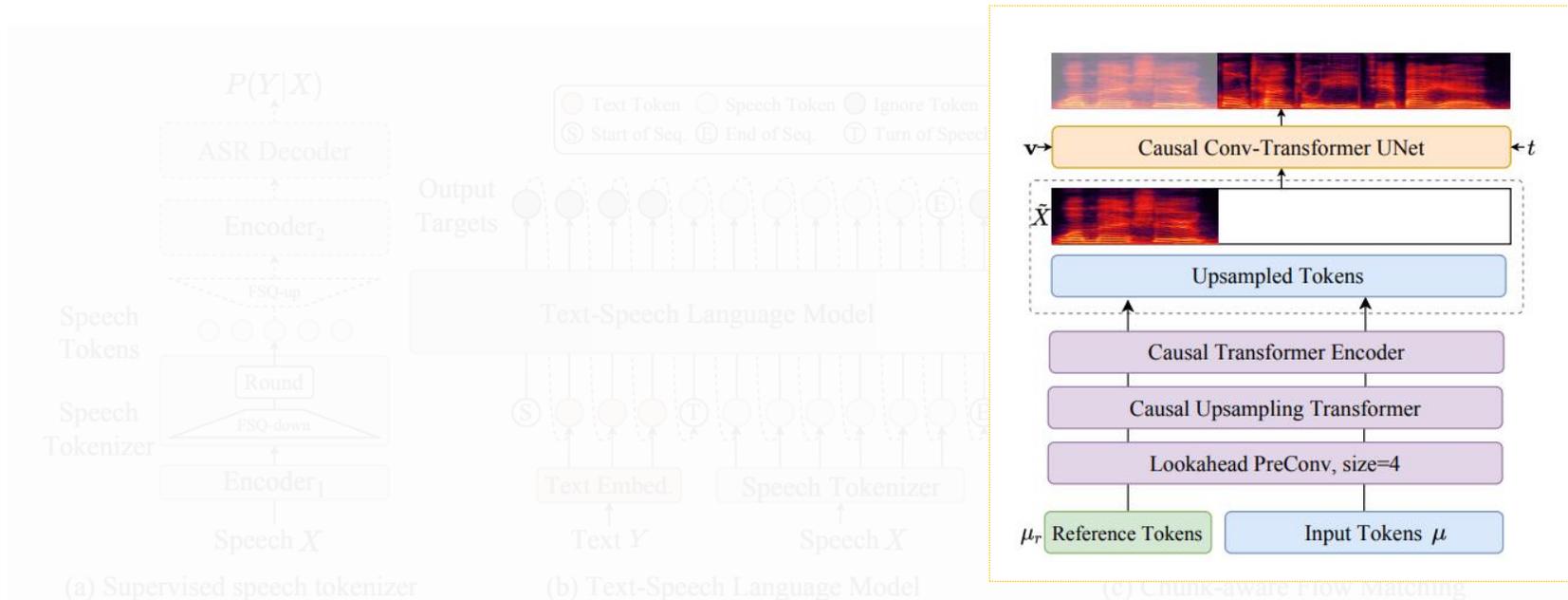
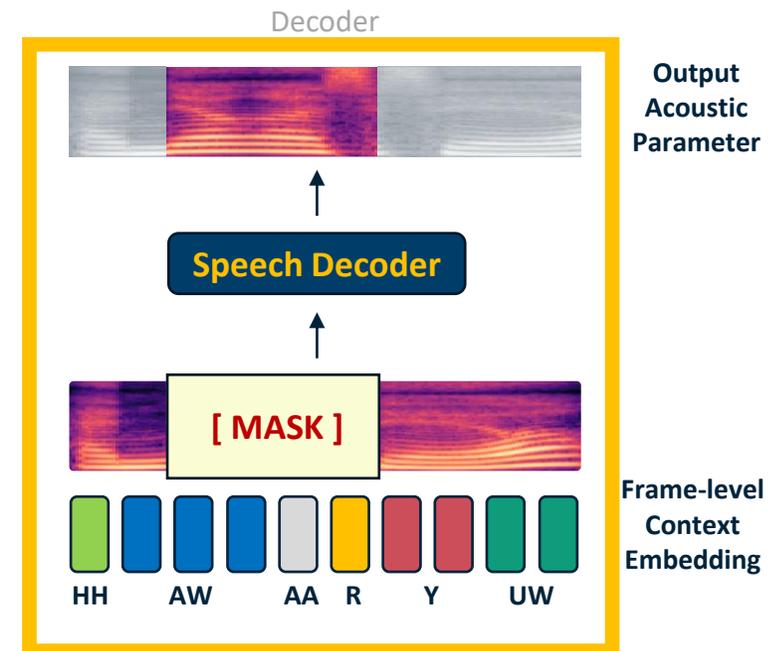
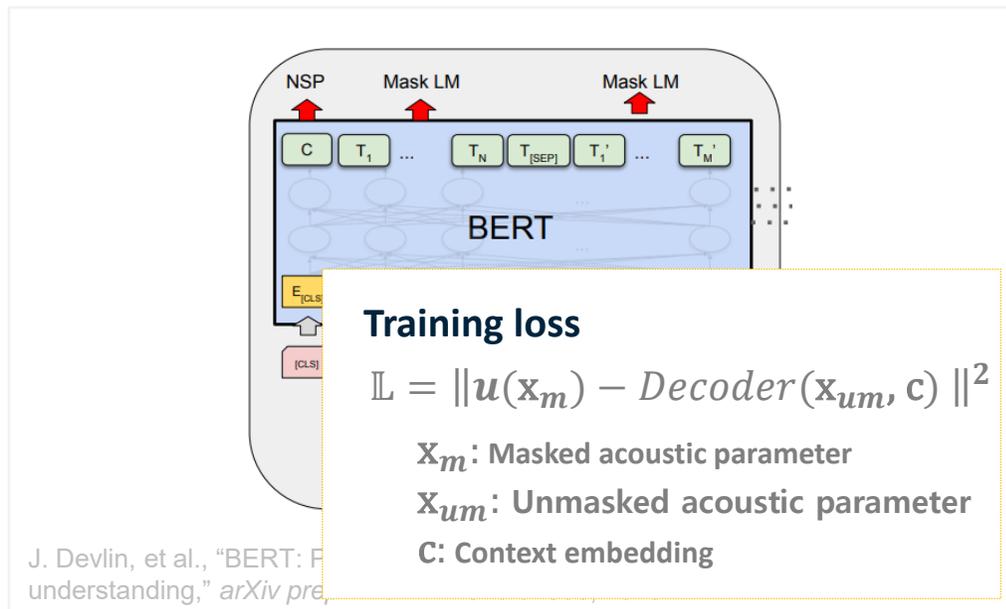


Figure 1: An overview of CosyVoice 2. (a) demonstrates the supervised speech tokenizer. Dashed modules are only used at the training stage. (b) is a unified text-speech language model for streaming and non-streaming synthesis. Dashed lines indicate the autoregressive decoding at the inference stage. (c) illustrates the causal flow matching model conditioning on a speaker embedding  $\mathbf{v}$ , semantic tokens  $\mu$ , masked speech features  $\tilde{X}$  and intermediate state  $X_t$  at timestep  $t$  on the probabilistic density path.

**Speech Token  $\rightarrow$  Speech**

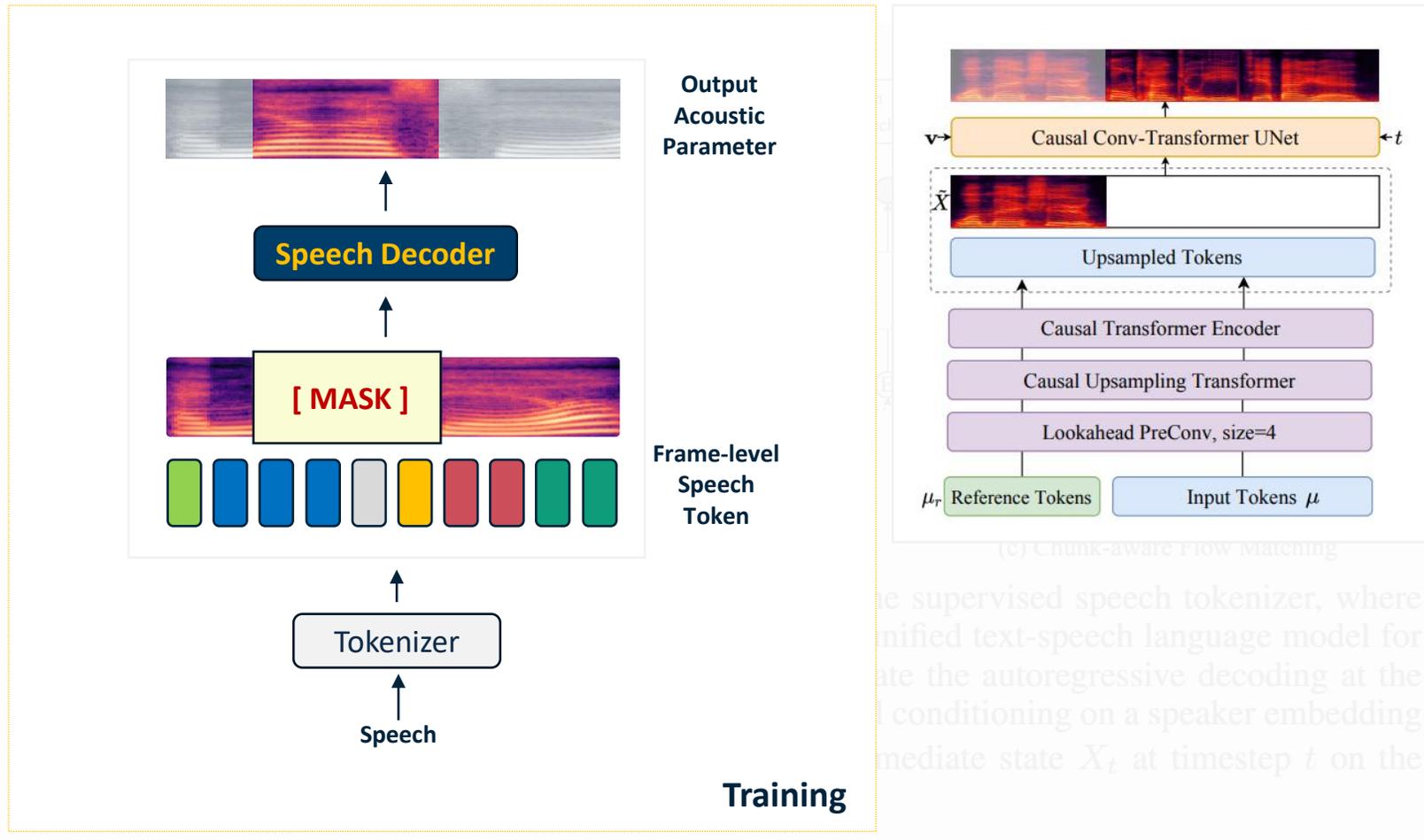
# Zero-shot voice cloning

Recall – Flow matching transformer



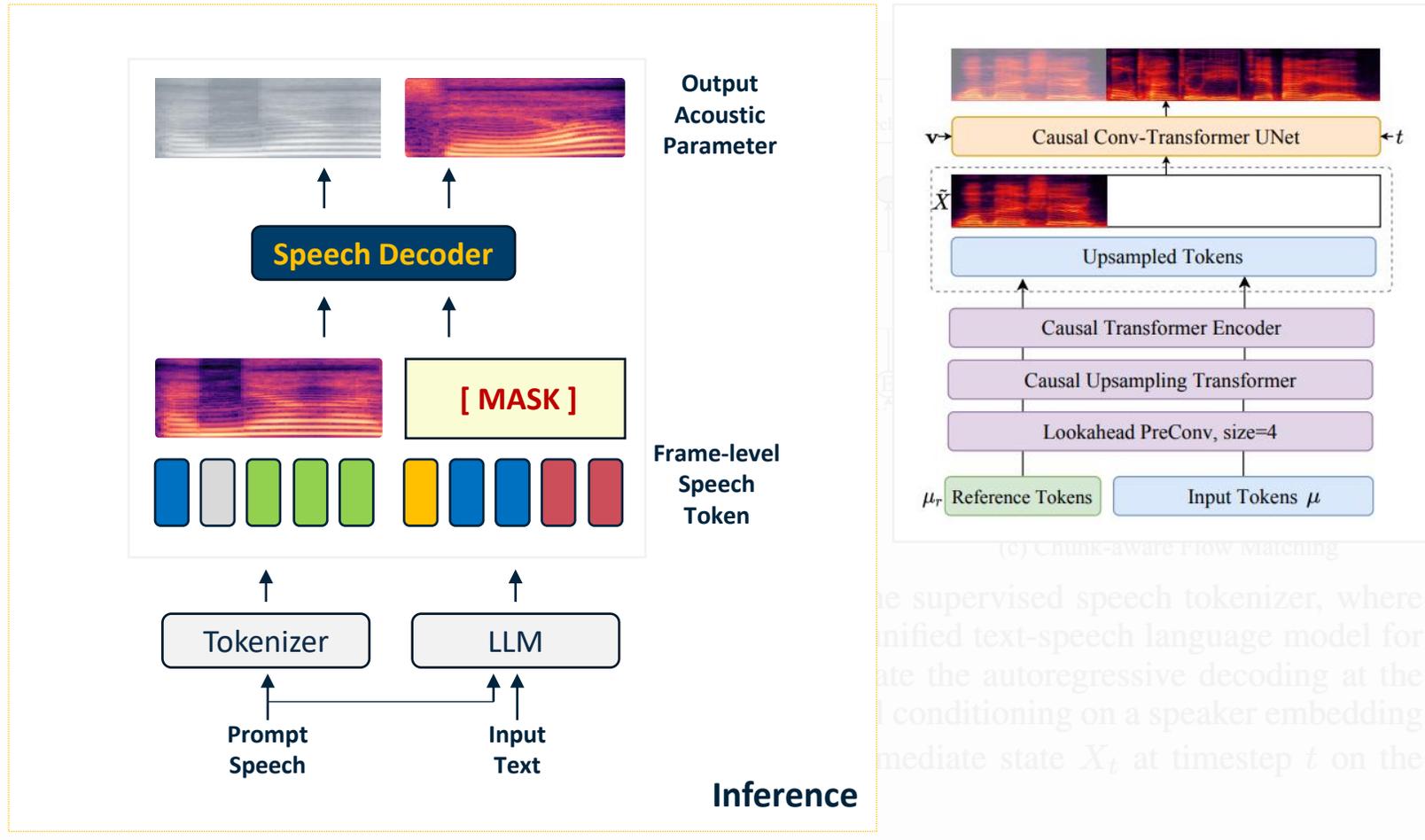
# Zero-shot voice cloning

## Speech decoder



# Zero-shot voice cloning

## Speech decoder



# Zero-shot voice cloning

## CosyVoice 2

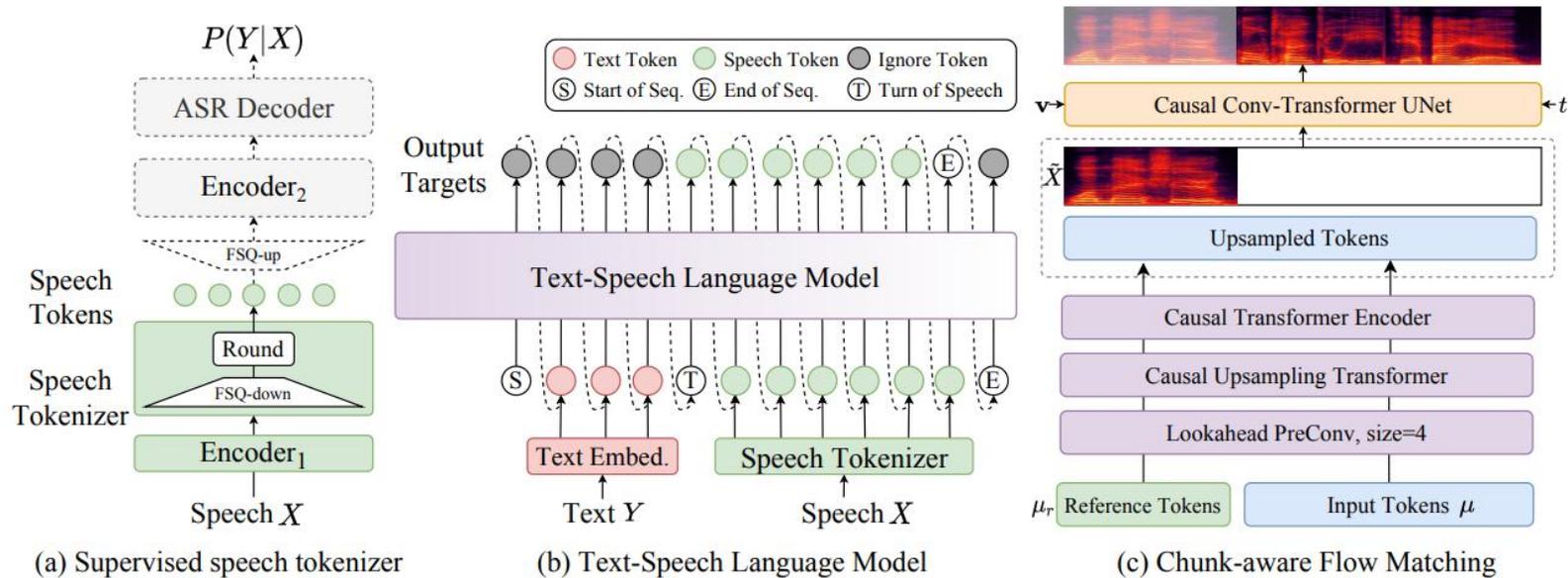


Figure 1: An overview of CosyVoice 2. (a) demonstrates the supervised speech tokenizer, where dashed modules are only used at the training stage. (b) is a unified text-speech language model for streaming and non-streaming synthesis. Dashed lines indicate the autoregressive decoding at the inference stage. (c) illustrates the causal flow matching model conditioning on a speaker embedding  $\mathbf{v}$ , semantic tokens  $\mu$ , masked speech features  $\tilde{X}$  and intermediate state  $X_t$  at timestep  $t$  on the probabilistic density path.

# Zero-shot voice cloning

## CosyVoice 2

Model	WER (%)	NMOS	SS
Human	2.66	3.84	0.697
ChatTTS [56]	6.84	3.89	-
GPT-SoVITs [57]	5.13	3.93	0.405
OpenVoice [58]	3.47	3.87	0.299
ParlerTTS [59]	3.16	3.86	-
EmotiVoice [60]	3.14	3.93	-
CosyVoice [34]	2.89	3.93	0.743
<b>CosyVoice 2</b>	<b>2.47</b>	<b>3.96</b>	<b>0.745</b>
<b>CosyVoice 2-S</b>	<b>2.45</b>	3.90	<b>0.751</b>

Table 5: Content consistency (WER), speaker similarity (SS) and speech quality (NMOS) results on LibriSpeech test-clean subset of baselines and CosyVoice 2. Whisper-Large V3 is employed as the ASR model and punctuations are excluded before WER calculation.

Model	<i>test-zh</i>		<i>test-en</i>		<i>test-hard</i>	
	CER (%) ↓	SS ↑	WER (%) ↓	SS ↑	WER (%) ↓	SS ↑
Human	1.26	0.755 (0.775)	2.14	0.734 (0.742)	-	-
Vocoder Resyn.	1.27	0.720	2.17	0.700	-	-
Seed-TTS <sup>†</sup> [33]	1.12	0.796	2.25	0.762	7.59	0.776
FireRedTTS [35]	1.51	0.635 (0.653)	3.82	0.460 (0.526)	17.45	0.621 (0.639)
MaskGCT [18]	2.27	0.774 (0.752)	2.62	0.714 (0.730)	10.27	0.748 (0.720)
E2 TTS (32 NFE) <sup>†</sup> [31]	1.97	0.730	2.19	0.710	-	-
F5-TTS (32 NFE) [32]	1.56	0.741 (0.794)	1.83	0.647 (0.742)	8.67	0.713 (0.762)
CosyVoice [34]	3.63	0.723 (0.775)	4.29	0.609 (0.699)	11.75	0.709 (0.755)
<b>CosyVoice 2</b>	<b>1.45</b>	<b>0.748 (0.806)</b>	<b>2.57</b>	<b>0.652 (0.736)</b>	<b>6.83</b>	<b>0.724 (0.776)</b>
<b>CosyVoice 2-S</b>	<b>1.45</b>	<b>0.753 (0.812)</b>	<b>2.38</b>	<b>0.654 (0.743)</b>	<b>8.08</b>	<b>0.732 (0.785)</b>

Table 6: Results of CosyVoice 2 and recent TTS models on the SEED test sets. † denotes close-sourced models. For speaker similarity, the result in a bracket are measured by ERes2Net, while the results outside brackets are measured by WavLM-based models.

# Zero-shot voice cloning

## CosyVoice 2

Strength

Fully enjoy LLM capability

Model	WER (%)	NMOS	SS
Human	2.66	3.84	0.697
ChatTTS [56]	6.84	3.89	-
GPT-SoVITS [57]	5.13	3.93	0.405
OpenVoice [58]	3.47	3.87	0.299
ParlerTTS [59]	3.16	3.86	-
EmotiVoice [60]	3.14	3.93	-
CosyVoice [34]	2.89	3.93	0.743
<b>CosyVoice 2</b>	<b>2.47</b>	<b>3.96</b>	<b>0.745</b>
CosyVoice 2-S	<b>2.45</b>	3.90	<b>0.751</b>

Table 5: Content consistency (WER), speaker similarity (SS) and speech quality (NMOS) results on LibriSpeech test-clean subset of baselines and CosyVoice 2. Whisper-Large V3 is employed as the ASR model and punctuations are excluded before WER calculation.

Model	<i>test-zh</i>		<i>test-en</i>		<i>test-hard</i>	
	CER (%) ↓	SS ↑	WER (%) ↓	SS ↑	WER (%) ↓	SS ↑
Human	1.26	0.755 (0.775)	2.14	0.734 (0.742)	-	-
Vocoder Resyn.	1.27	0.720	2.17	0.700	-	-
Seed-TTS <sup>†</sup> [33]	1.12	0.796	2.25	0.762	7.59	0.776
FireRedTTS [35]	1.51	0.635 (0.653)	3.82	0.460 (0.526)	17.45	0.621 (0.639)
MaskGCT [18]	2.27	0.774 (0.752)	2.62	0.714 (0.730)	10.27	0.748 (0.720)
E2 TTS (32 NFE) <sup>†</sup> [31]	1.97	0.730	2.19	0.710	-	-
F5-TTS (32 NFE) [32]	1.56	0.741 (0.794)	1.83	0.647 (0.742)	8.67	0.713 (0.762)
CosyVoice [34]	3.63	0.723 (0.775)	4.29	0.609 (0.699)	11.75	0.709 (0.755)
<b>CosyVoice 2</b>	<b>1.45</b>	<b>0.748 (0.806)</b>	<b>2.57</b>	<b>0.652 (0.736)</b>	<b>6.83</b>	<b>0.724 (0.776)</b>
CosyVoice 2-S	1.45	0.753 (0.812)	2.38	0.654 (0.743)	8.08	0.732 (0.785)

Table 6: Results of CosyVoice 2 and recent TTS models on the SEED test sets. † denotes close-sourced models. For speaker similarity, the result in a bracket are measured by ERes2Net, while the results outside brackets are measured by WavLM-based models.

# Zero-shot voice cloning

CosyVoice 2

## Strength

Fully enjoy LLM capability

## Weakness

Hallucinations  
from the model's autoregressive nature

Flow matching  
(Non-AR) 😊



CosyVoice 2  
(AR) 😞



# Zero-shot voice cloning

CosyVoice 2

## Strength

Fully enjoy LLM capability

## Weakness

Hallucinations  
from the model's autoregressive nature

Flow matching  
(Non-AR) 😊



CosyVoice 2  
(AR) 😞



# Zero-shot voice cloning

CosyVoice 2

## Strength

Fully enjoy LLM capability

## Weakness

Hallucinations  
from the model's autoregressive nature

## Solution?

Data refinement  
Reinforcement learning

...

# Zero-shot voice cloning

CosyVoice 2

## Strength

Fully enjoy LLM capability

CosyVoice 2  
(Official) 😞

CosyVoice 2  
(Refined) 😊



## Solution?

Data refinement

Reinforcement learning

...

# Zero-shot voice cloning

CosyVoice 2

## Strength

Fully enjoy LLM capability

CosyVoice 2  
(Official) 😞



CosyVoice 2  
(Refined) 😊



## Solution?

Data refinement

Reinforcement learning

...

# Zero-shot voice cloning

CosyVoice 2

## Strength

Fully enjoy LLM capability

CosyVoice 2  
(Official) 😞

CosyVoice 2  
(Refined) 😊



## Solution?

Data refinement

Reinforcement learning

...

# Zero-shot voice cloning

CosyVoice 2

## Strength

Fully enjoy LLM capability

CosyVoice 2  
(Official) 😞



CosyVoice 2  
(Refined) 😊



## Solution?

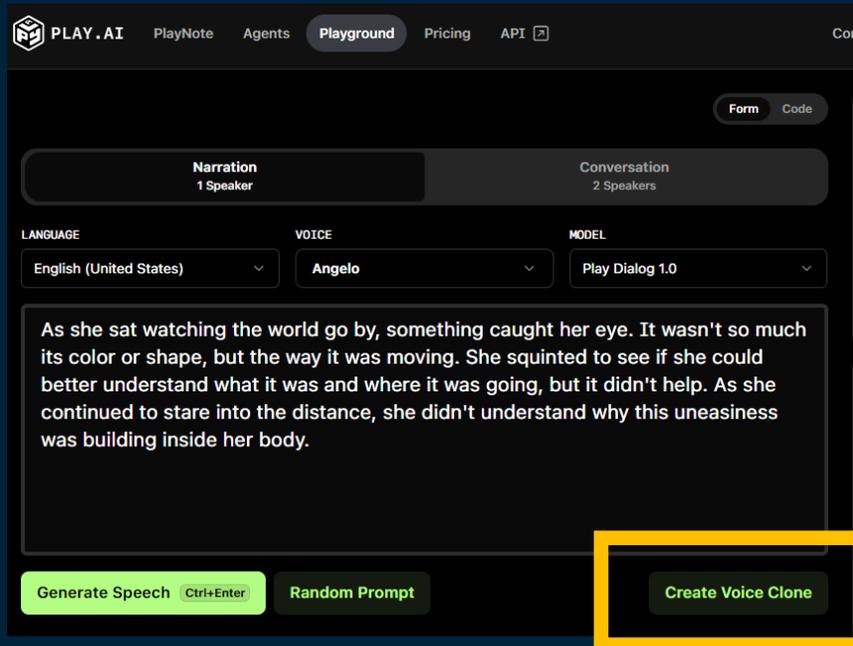
Data refinement

Reinforcement learning

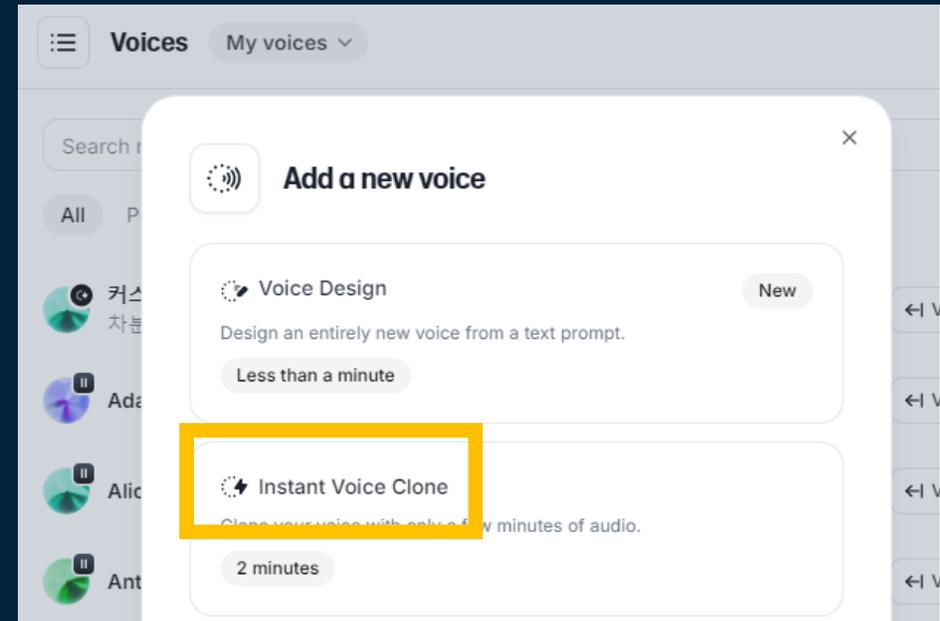
...

# Scaling TTS model

## Voice cloning



<https://play.ht/>



<https://elevenlabs.io/>

# Scaling TTS model

## Voice cloning



The screenshot shows the Play.ht website interface. At the top, there is a navigation bar with 'PLAY.ht' and 'My voices'. Below this, there are several tabs: 'Voice', 'Narration', 'Text-to-speech', and 'Text-to-video'. Under the 'Voice' tab, there are three voice options: 'English (American)', 'Angela', and 'English (British)'. Below the voice options, there is a text input field with the text: 'As she sat watching the world go by, something caught her eye. It wasn't so much its color or shape, but the way it was moving. She squinted to see if she could better understand what it was and where it was going, but it didn't help. As she continued to stare into the distance, she didn't understand why this uneasiness was building inside her body.' At the bottom, there is a 'Generate sample audio' button and a 'Random prompt' button.

<https://play.ht/>



The screenshot shows the ElevenLabs website interface. At the top, there is a navigation bar with 'Voices' and 'My voices'. Below this, there is a large 'Add a new voice' button. Underneath, there are two options: 'Voice Design' and 'Instant Voice Clone'. The 'Voice Design' option is described as 'Design an entirely new voice from a text prompt' and 'Less than a minute'. The 'Instant Voice Clone' option is described as 'Clone a voice from an audio sample' and 'Less than 5 minutes of audio'. At the bottom, there is a 'Generate sample audio' button.

<https://elevenlabs.io/>

**Ethical problem ?**

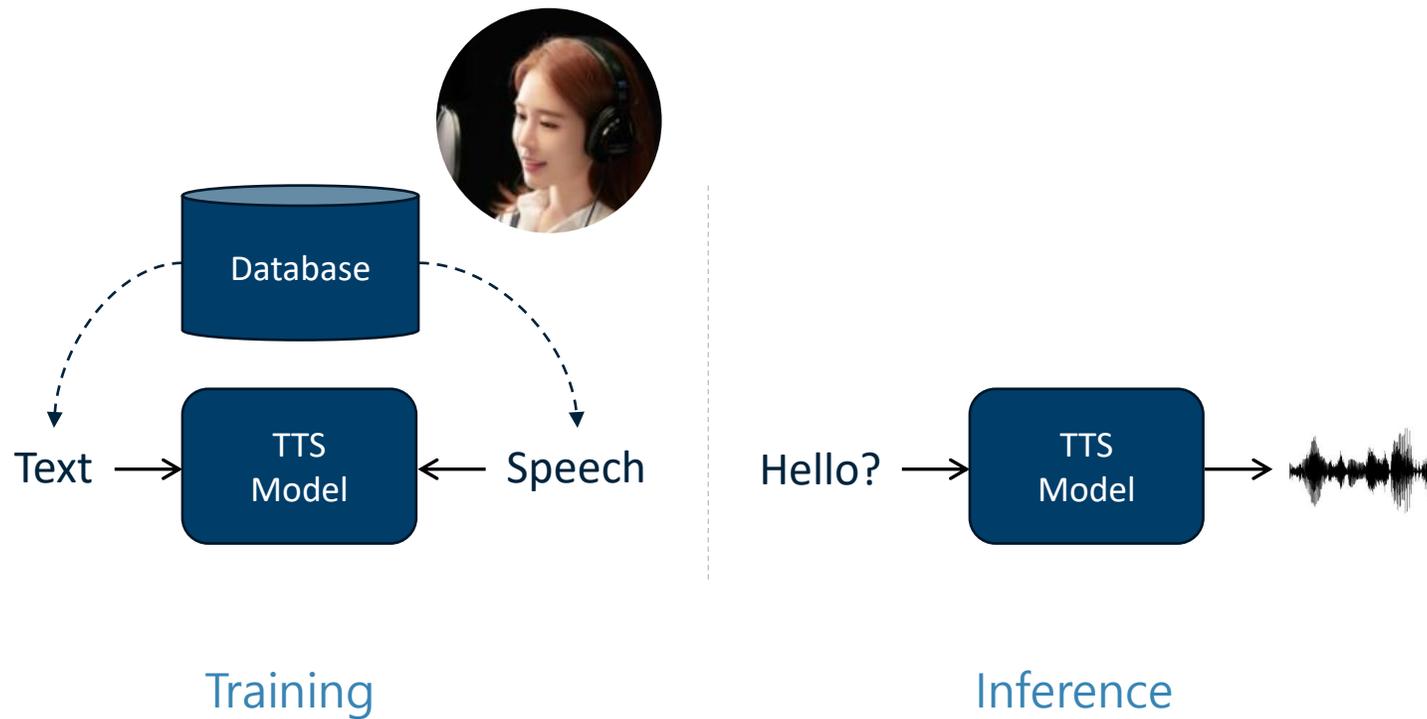


# Summary

---

# Summary

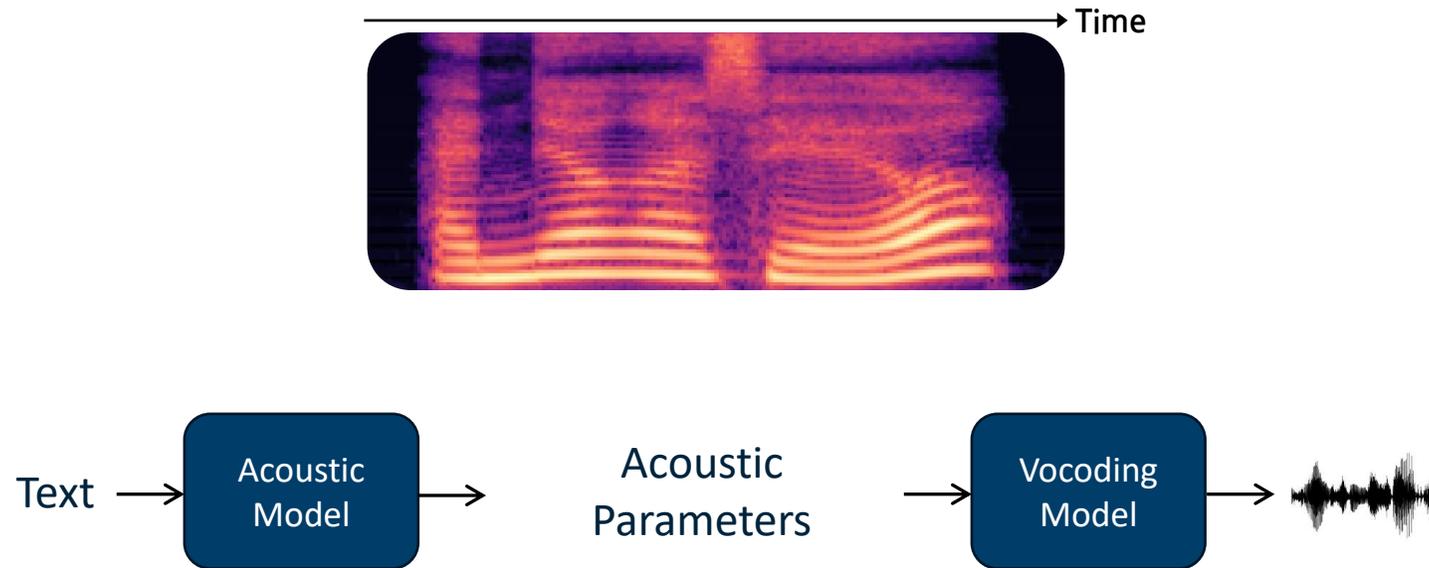
## Deep learning-based TTS system



**Human-like voice quality** 😊

# Summary

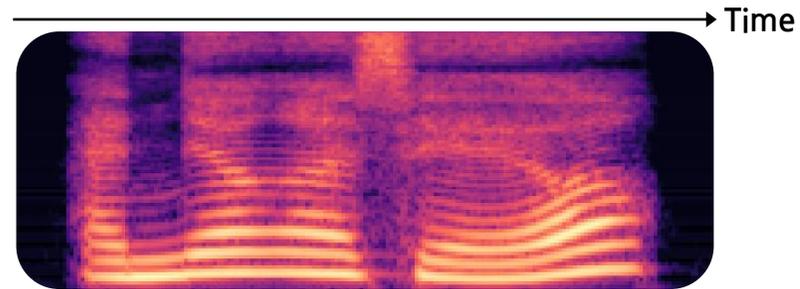
## Deep learning-based TTS system



**Acoustic model + Vocoding model**

# Summary

## Deep learning-based TTS system



Estimating acoustic parameters from text inputs

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

Acoustic model + Vocoding model

# Summary

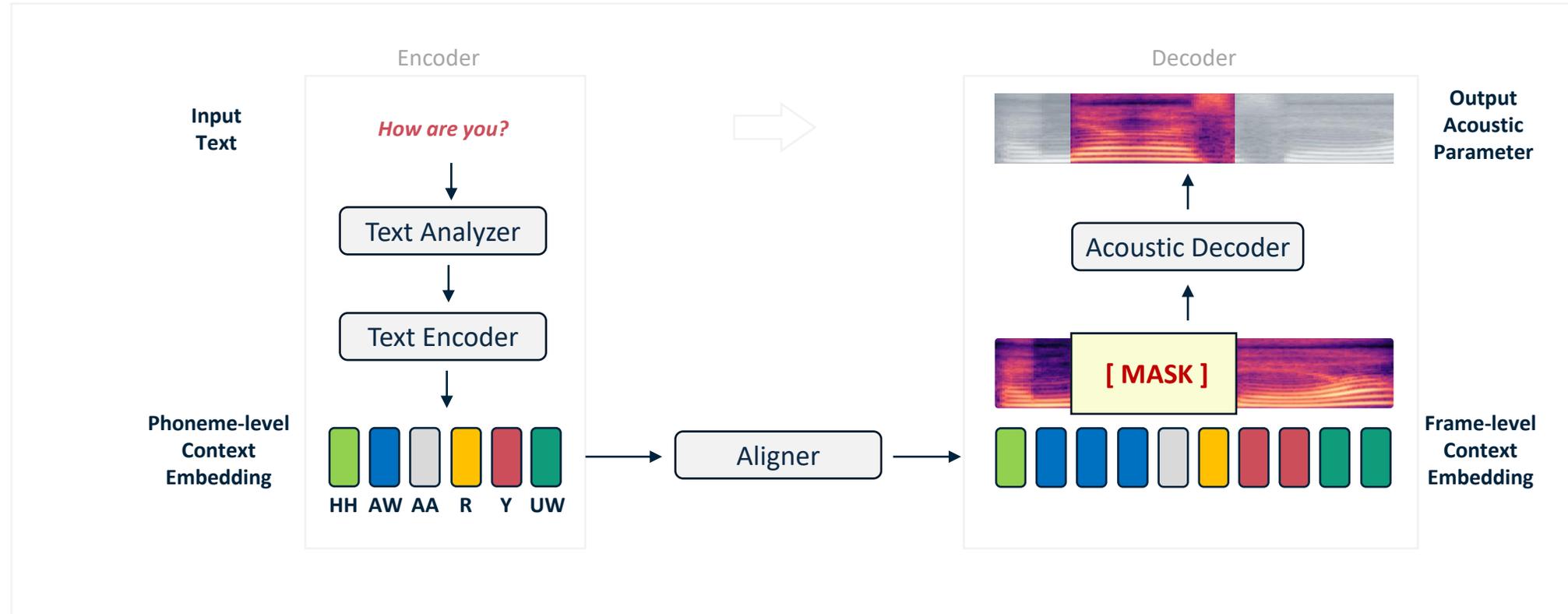
## Recording constraint

	Conventional TTS	Voice cloning
Recording amount	> 30~60 min	< Few seconds
Speaking type	Script reading	Spontaneous speaking
Speaker	Professional voice actor	Non-professional
Recording amount	Clean studio	Anywhere
TTS quality	Natural	<b>Very natural</b>

~~Recording quality matters: Poor recording → TTS degradation~~

# Summary

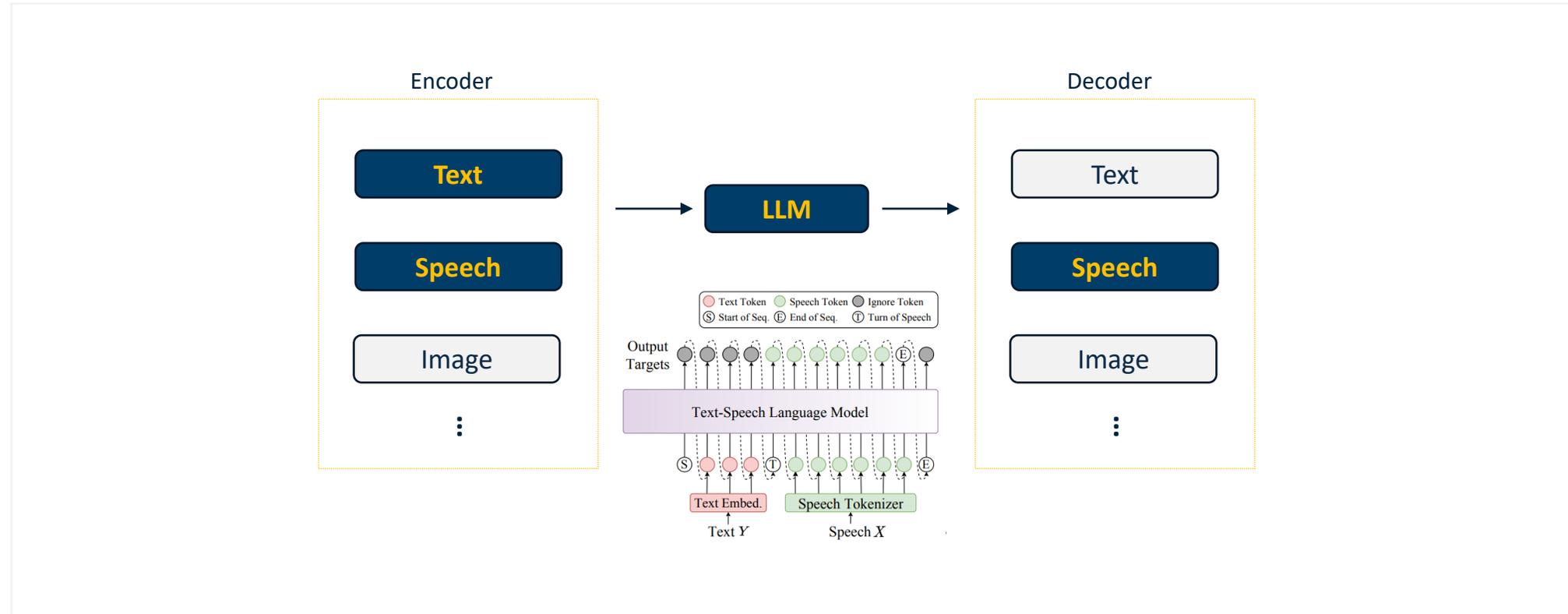
## Applying audio infilling task



The model focuses on **relationship** between **adjacent acoustic parameters**, rather than reconstructing the target data

# Summary

## Applying LLM task



The model fully enjoys the **benefit of LLM capability**

Q / A



[gregorio.song@gmail.com](mailto:gregorio.song@gmail.com)

