


# The Sound Dimension: Speech and Audio in Multimodal AI

Vineet Gandhi

Associate Professor

Center for Visual Information Technology





Why?  
How?  
What?  
Where?

Speech and Audio in MAI

# Excellent survey paper

*Baltrusaitis et al.* Multimodal Machine Learning: A Survey and Taxonomy. **TPAMI 2019**







# VOICE OR TEXT?

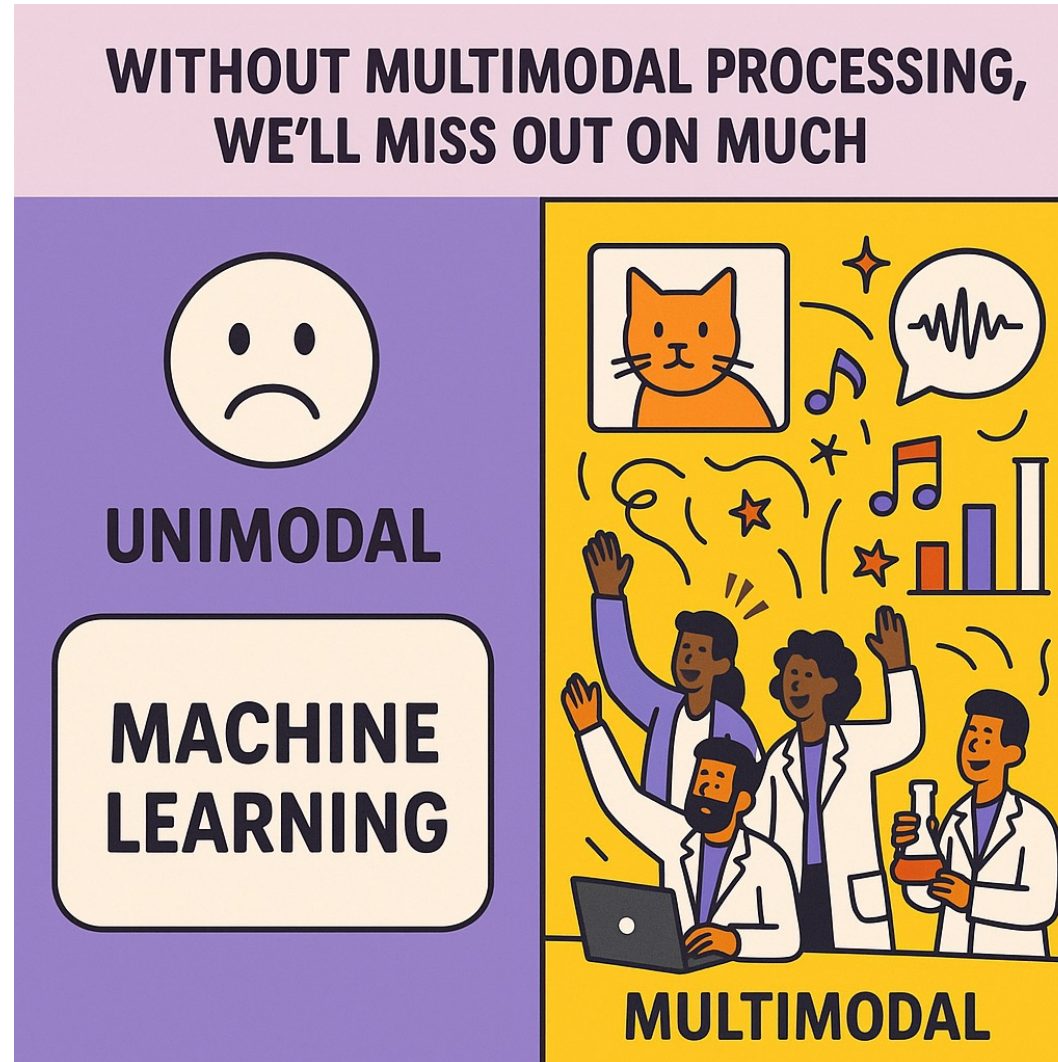




# HUMAN LEARNING IS INHERENTLY MULTIMODAL

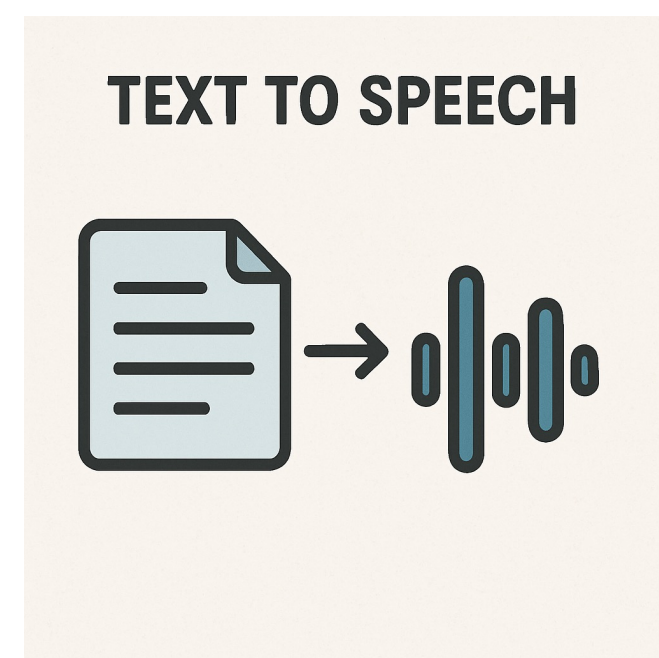
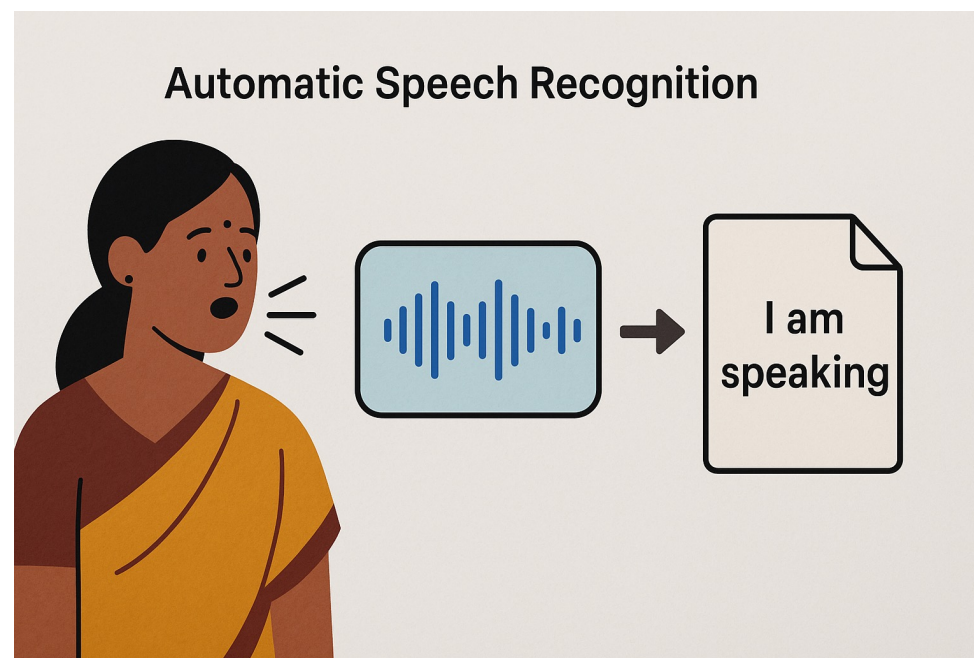


And more importantly, we will miss out a lot of fun applications/use-cases





# Key Speech problems

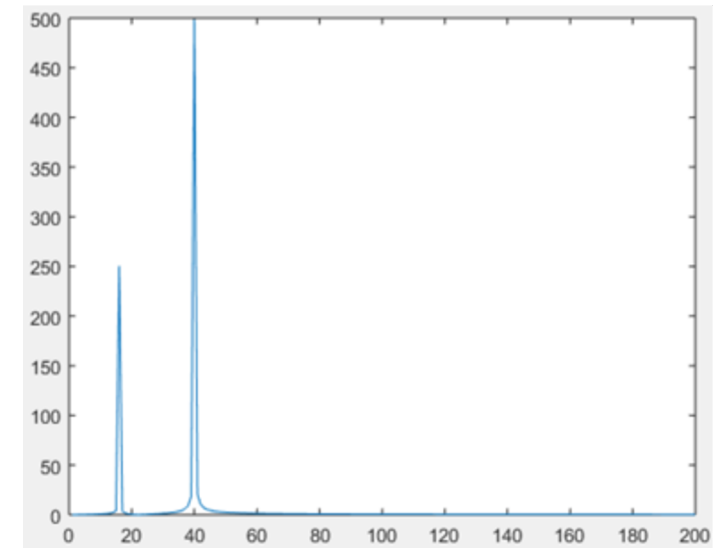
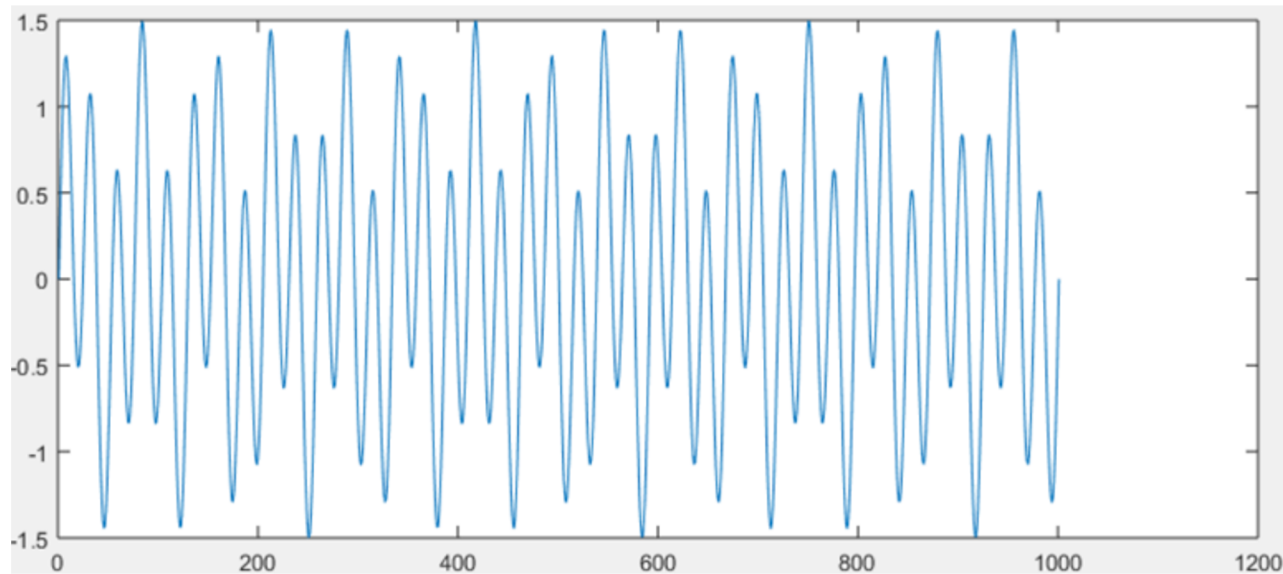


# Multimodal processing has transformed Speaker Recognition



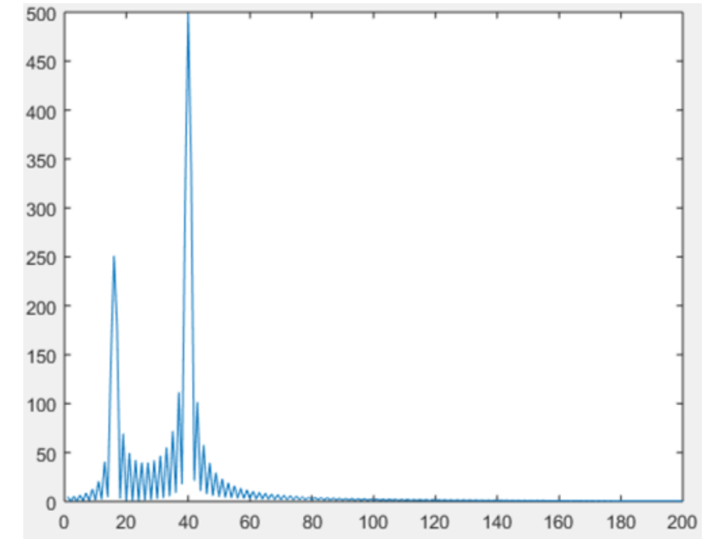
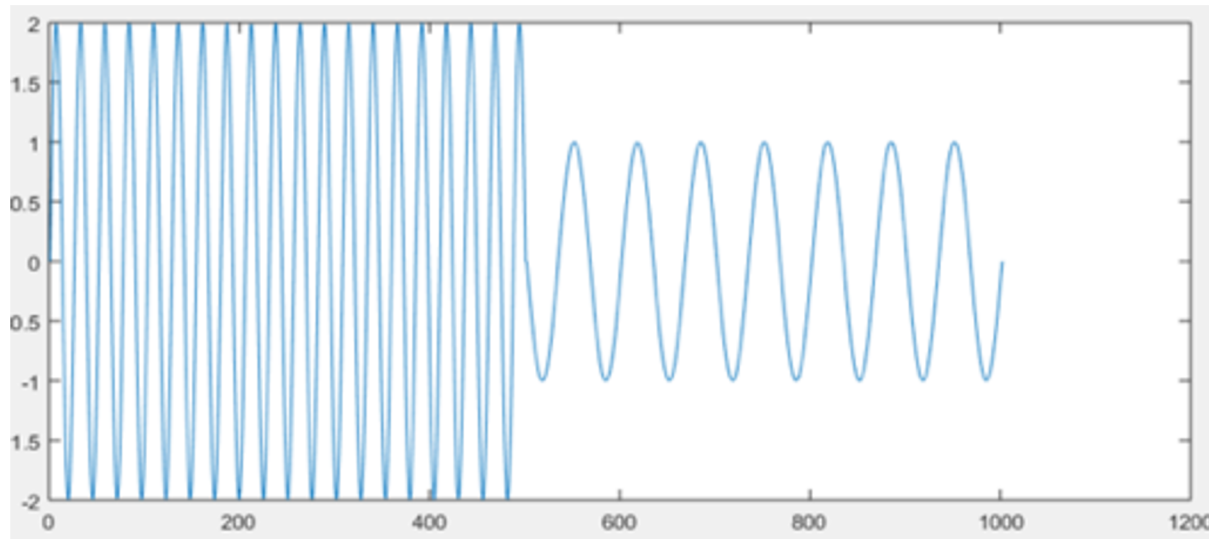
# Some basics to help my argument

$$f(t) = \sin(2\pi \cdot 39t) + 0.5 \sin(2\pi \cdot 15t)$$



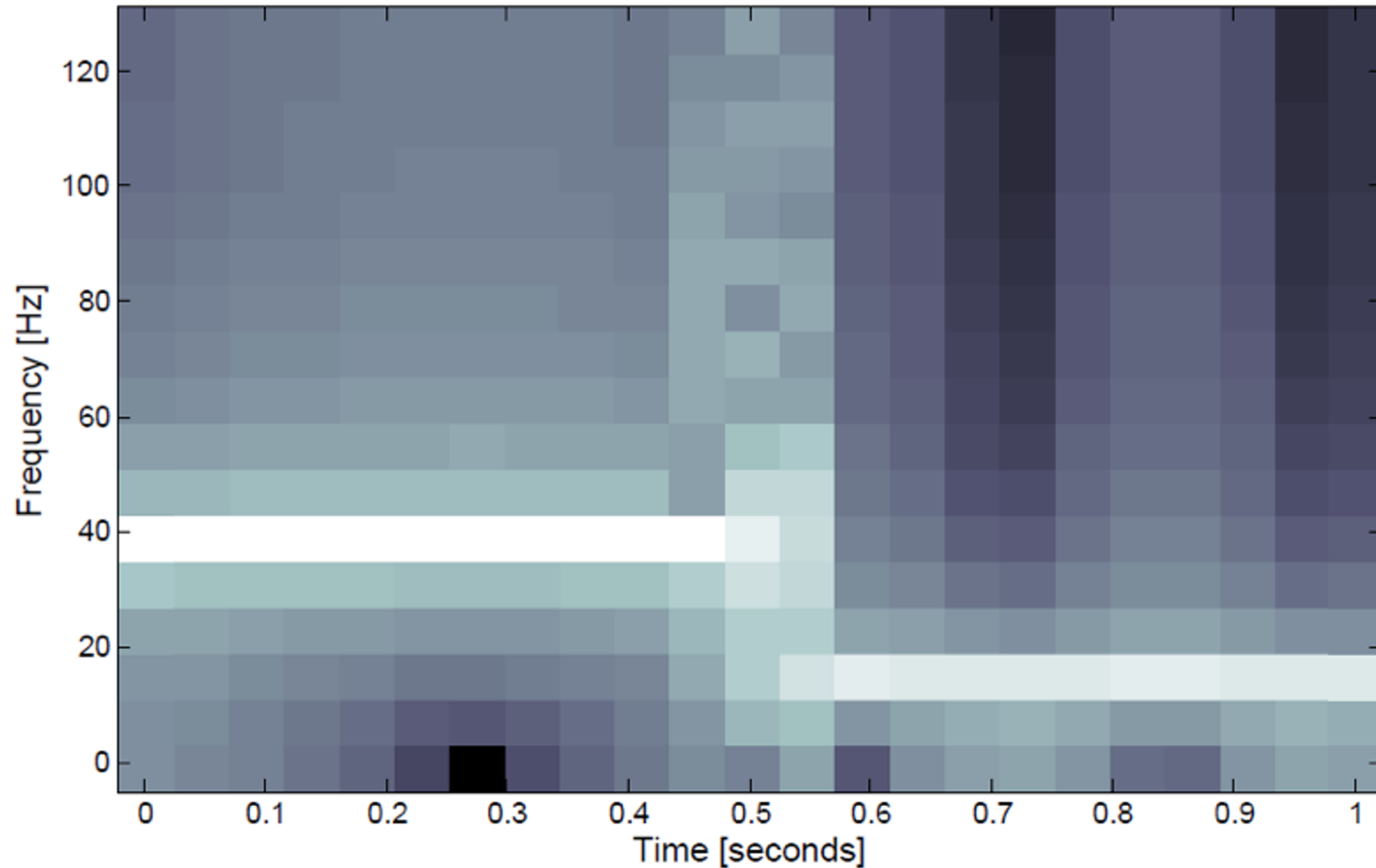
# Another example sound signal

$$g(t) = \begin{cases} 2 * \sin(2\pi \cdot 39t), & 0 \leq t \leq 1/2 \\ \sin(2\pi \cdot 15t), & 1/2 < t \leq 1 \end{cases}$$





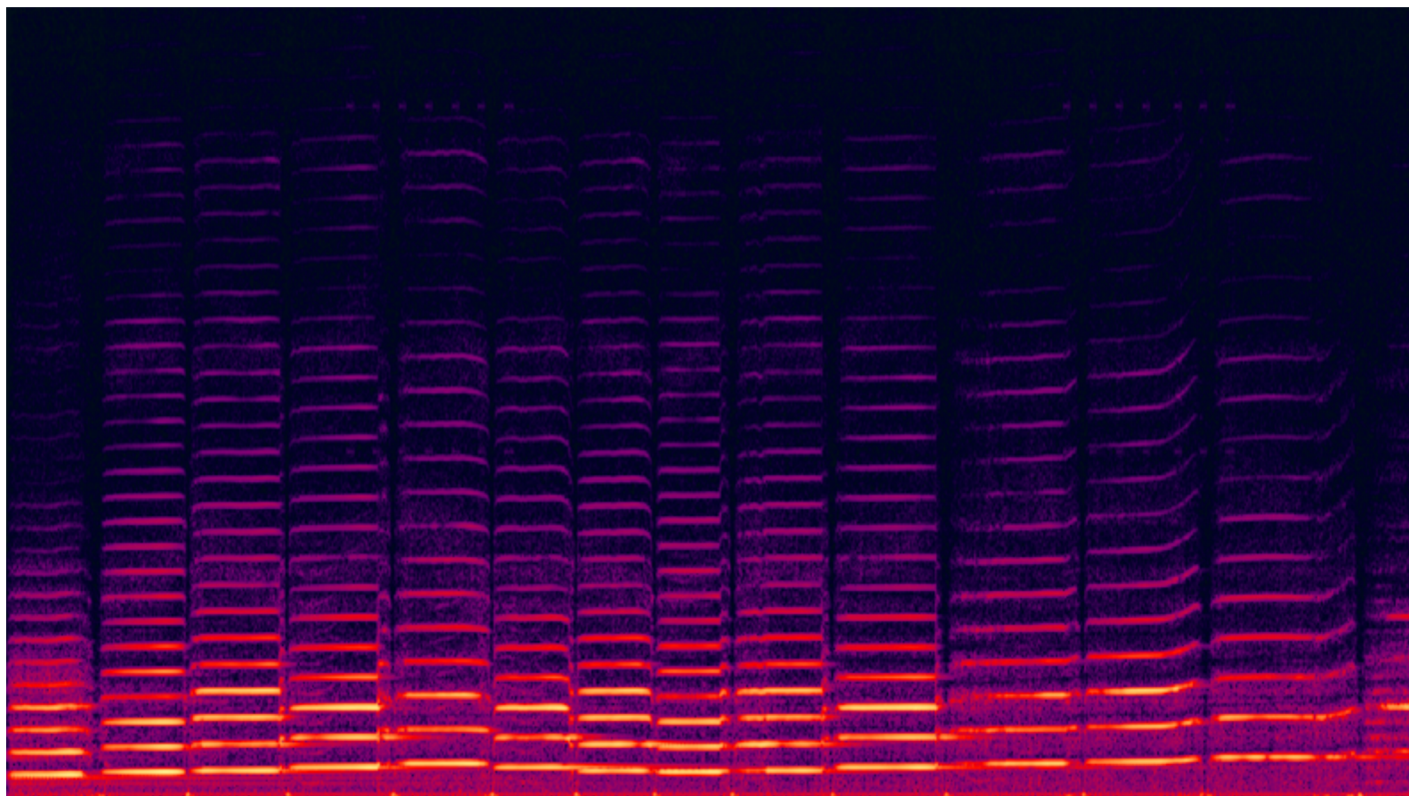
# Spectrogram



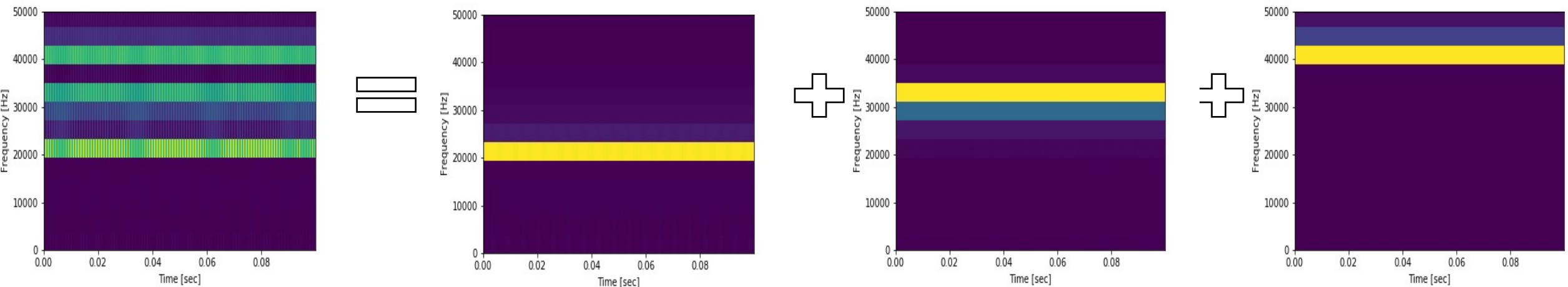
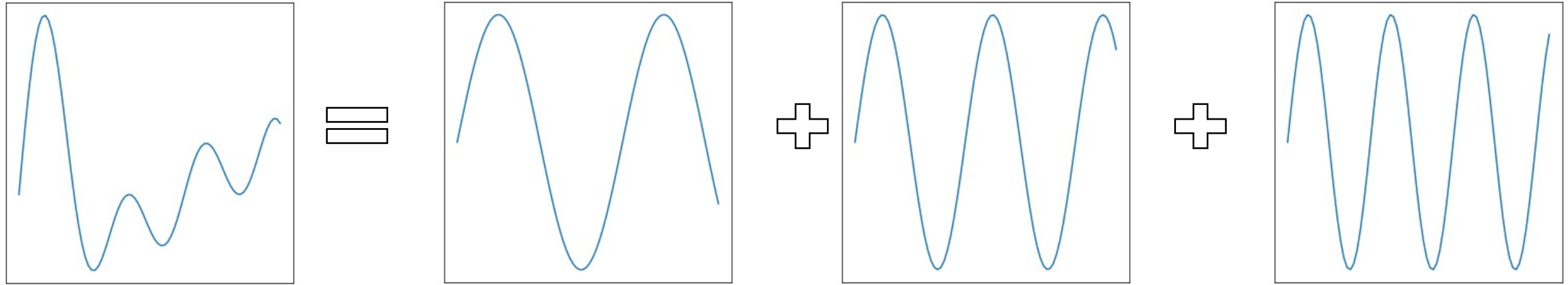
Spectrogram of a piecewise monochromatic signal.

Lighter color indicates greater DFT magnitude

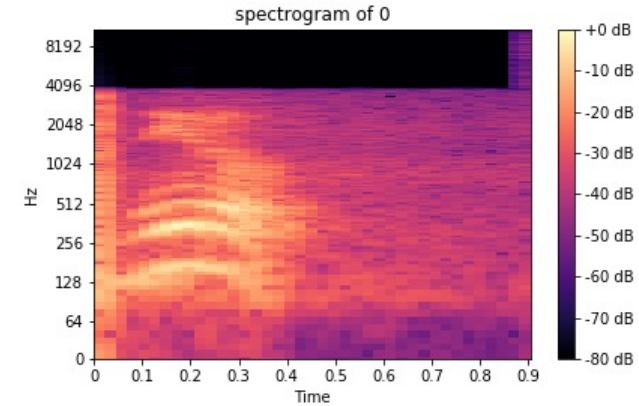
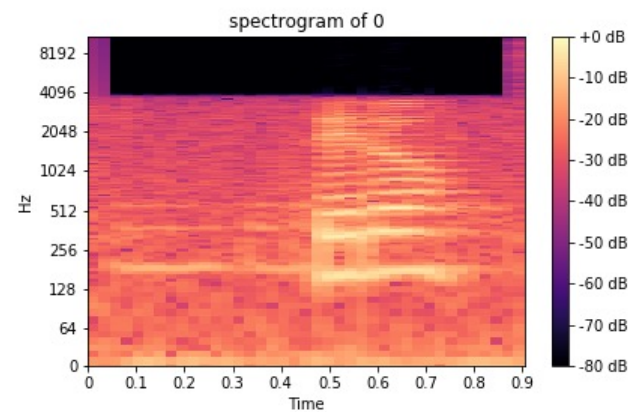
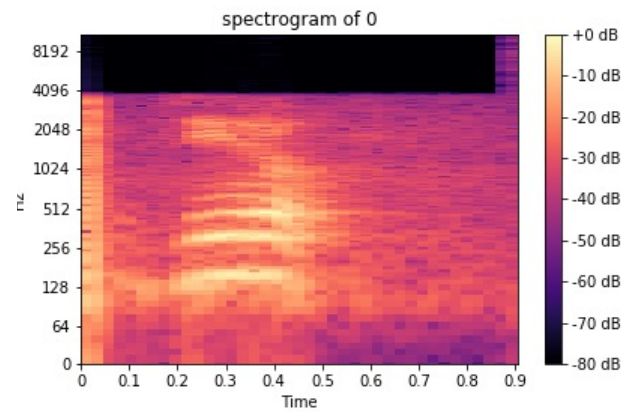
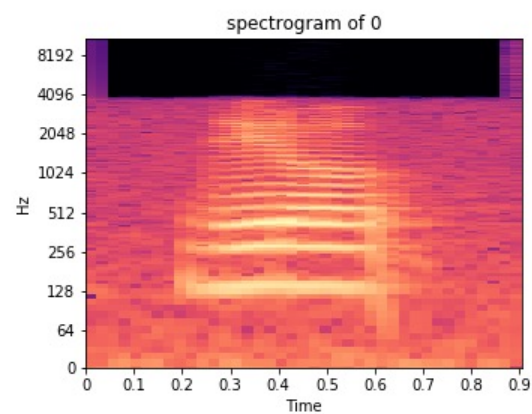
# Spectrogram



# Wave as a combination of sine waves

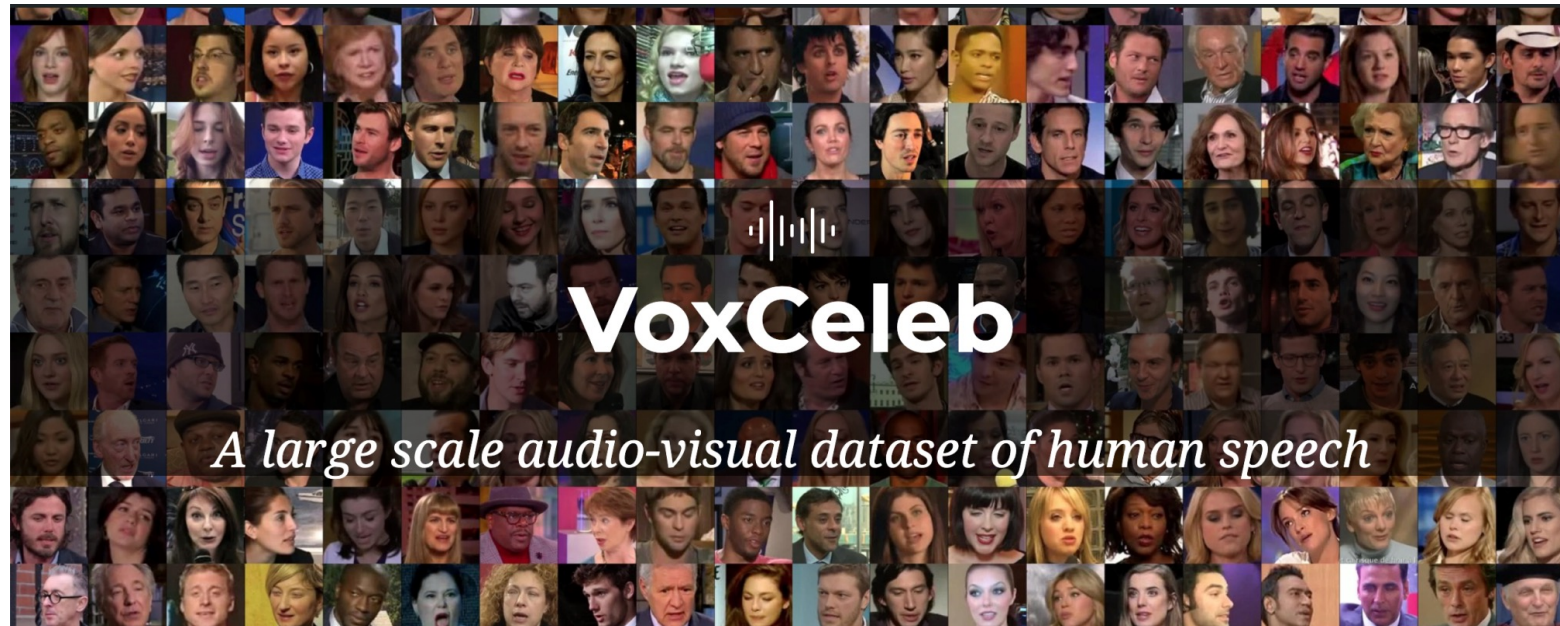


# Utterance of word zero



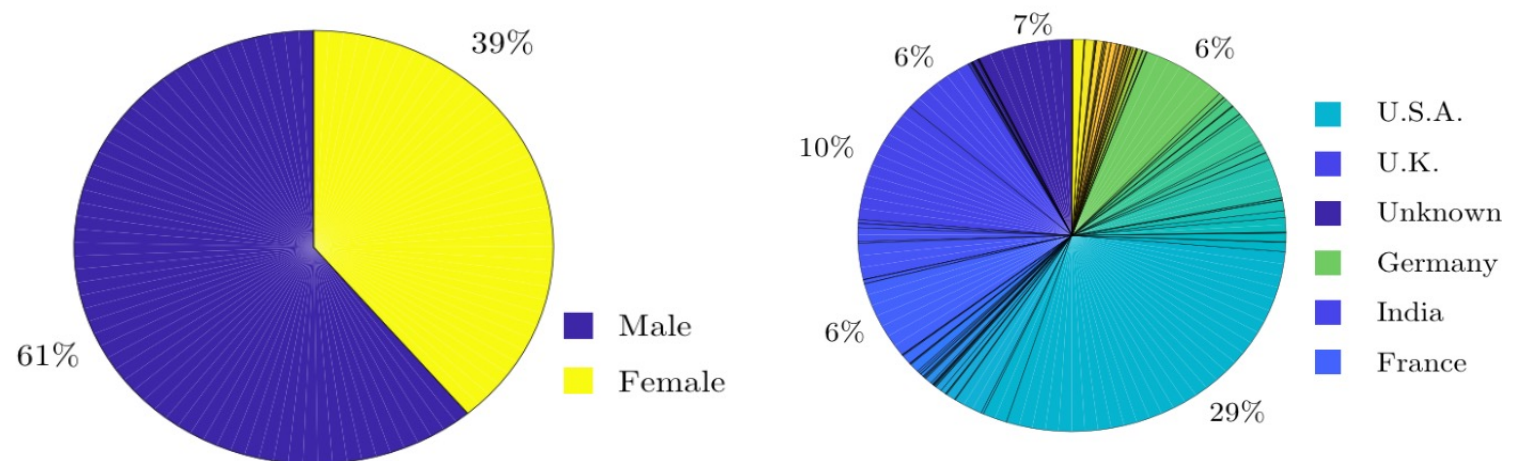


# Transforming speaker verification/identification

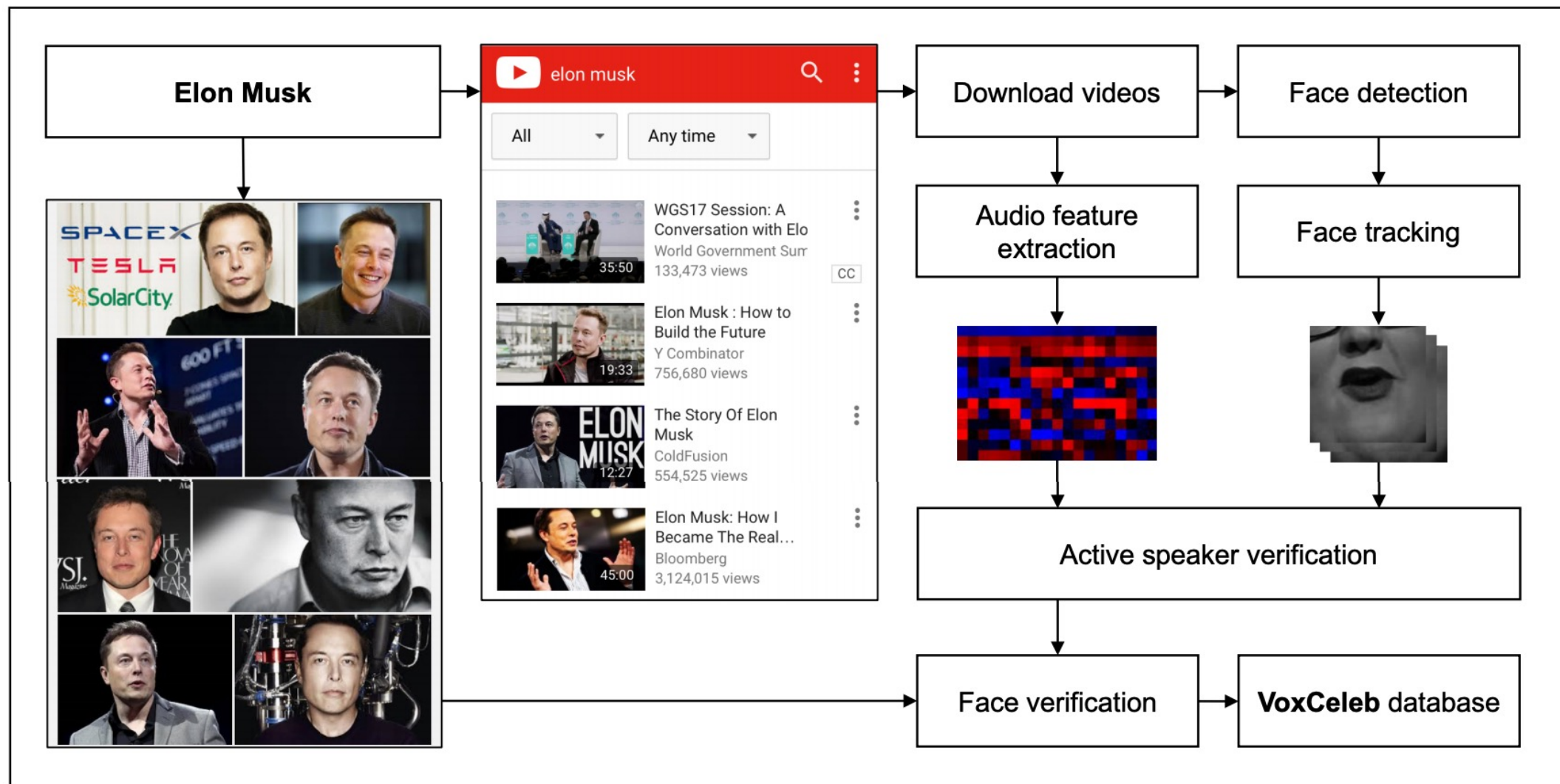


## VoxCeleb2

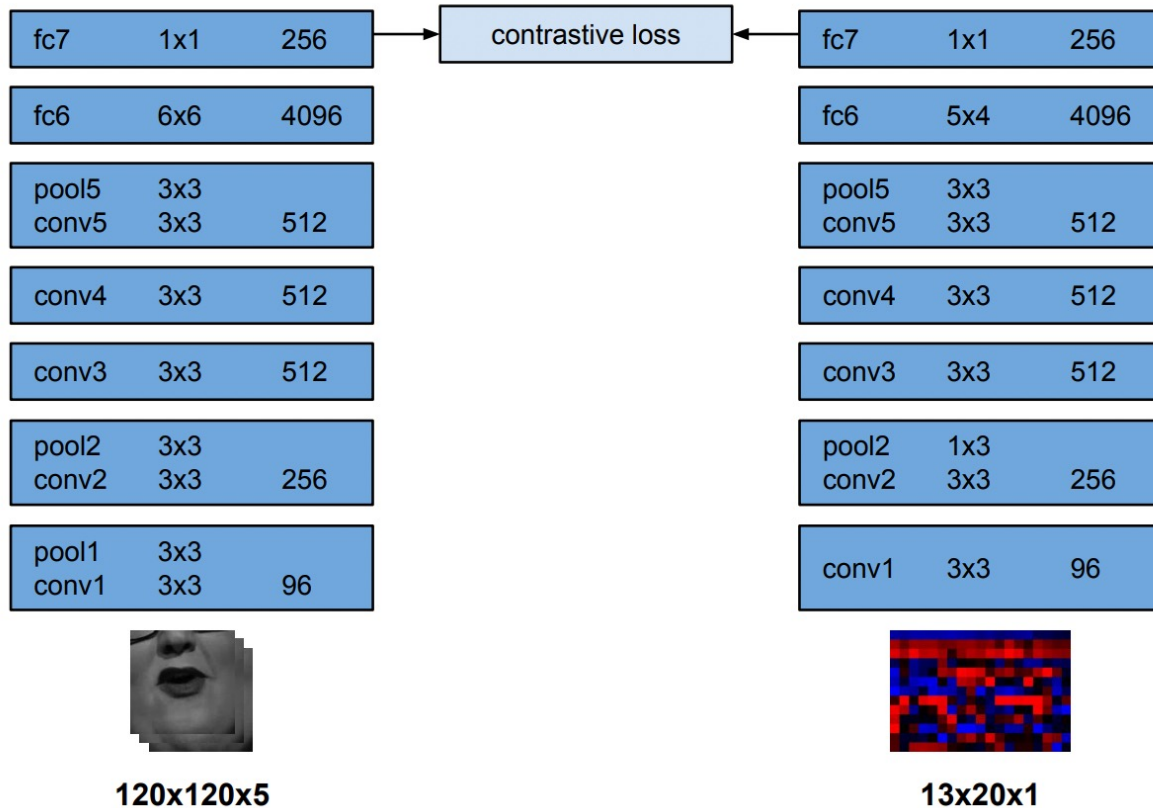
*VoxCeleb2 contains over a million utterances for 6,112 identities.*



# VoxCeleb: automated data collection



# SyncNet



## Solves three tasks

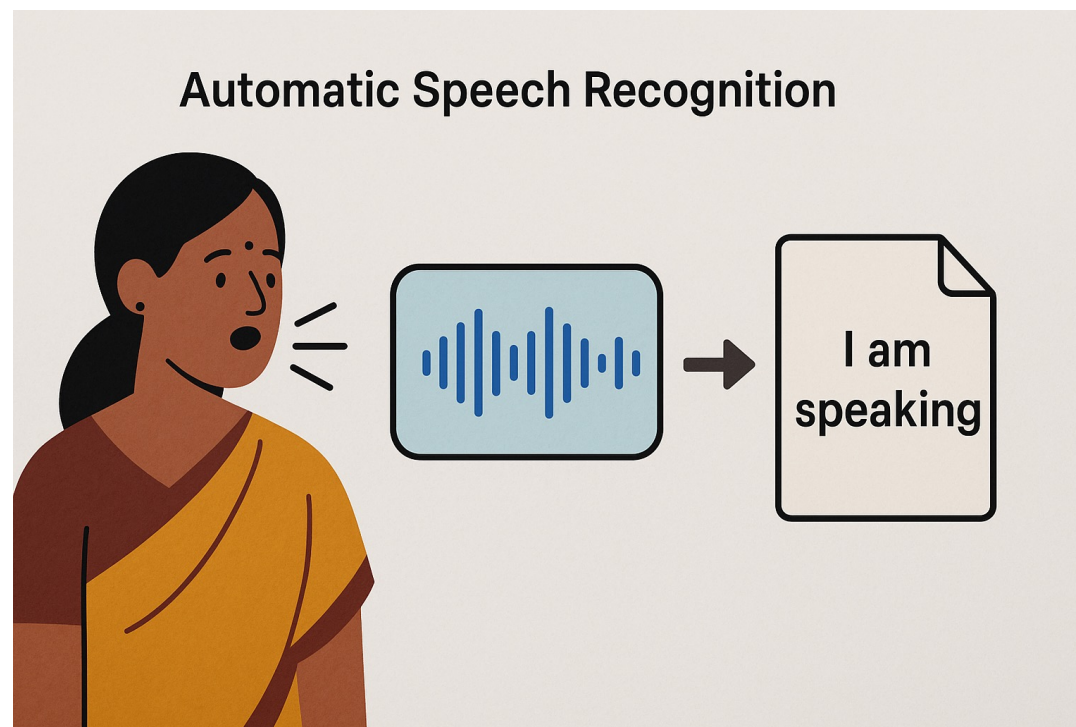
- Determining the lip-sync error in videos
- Detecting the speaker in a scene with multiple faces
- Lip reading

# VoxCeleb performance

<b>Accuracy</b>	<b>Top-1 (%)</b>	<b>Top-5 (%)</b>
<b>I-vectors + SVM</b>	49.0	56.6
<b>I-vectors + PLDA + SVM</b>	60.8	75.6
<b>CNN-fc-3s no var. norm.</b>	63.5	80.3
<b>CNN-fc-3s</b>	72.4	87.4
<b>CNN</b>	<b>80.5</b>	<b>92.1</b>



# Transformation through transformers



# Transformers and Self-Supervision

1. Shared pool of architectural insights
  - Models like **HuBERT**, **wav2vec 2.0**, and **Whisper** borrow architectural insights (e.g., Transformers, masked modeling) from NLP
2. Language Models Enhance ASR Decoding
  - Leads to **better handling of rare words, disfluencies, and long-range dependencies**

# HuBERT Training Process

Alternate between two steps

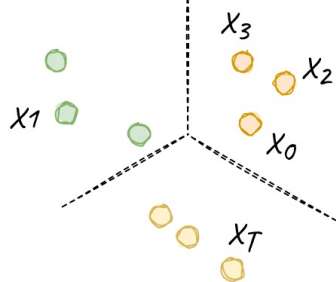
## STEP 1: Discover "hidden units" targets

Hidden units embeddings

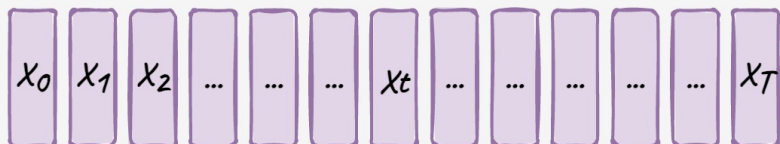


Use hidden units as targets to predict

### K-Means Clustering



Assign each feature vector  $X_t$  to a hidden unit cluster



### Clustering Feature Extraction

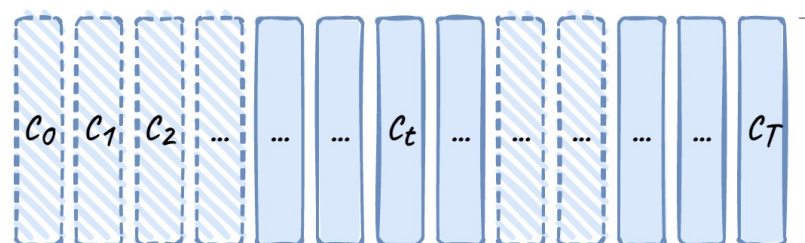
Audio Waveform



Compute directly from raw waveform

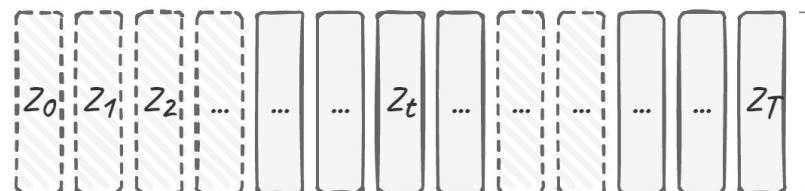
## STEP 2: Predict targets at masked positions

### Cross-Entropy Loss (Predict hidden units at masked locations)



Use context representations for prediction

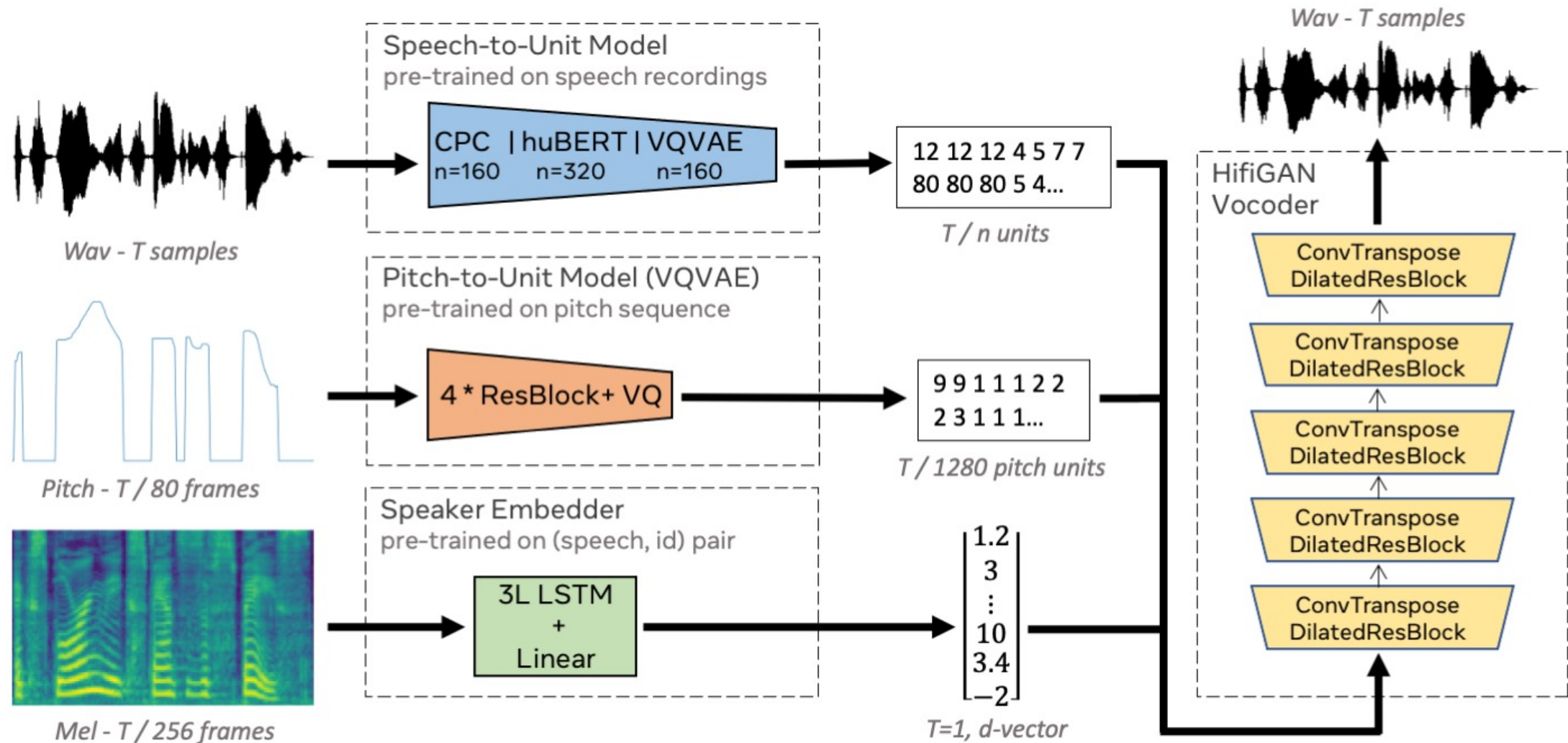
### Context Network (Transformer Encoder)



Mask ~50% of time steps

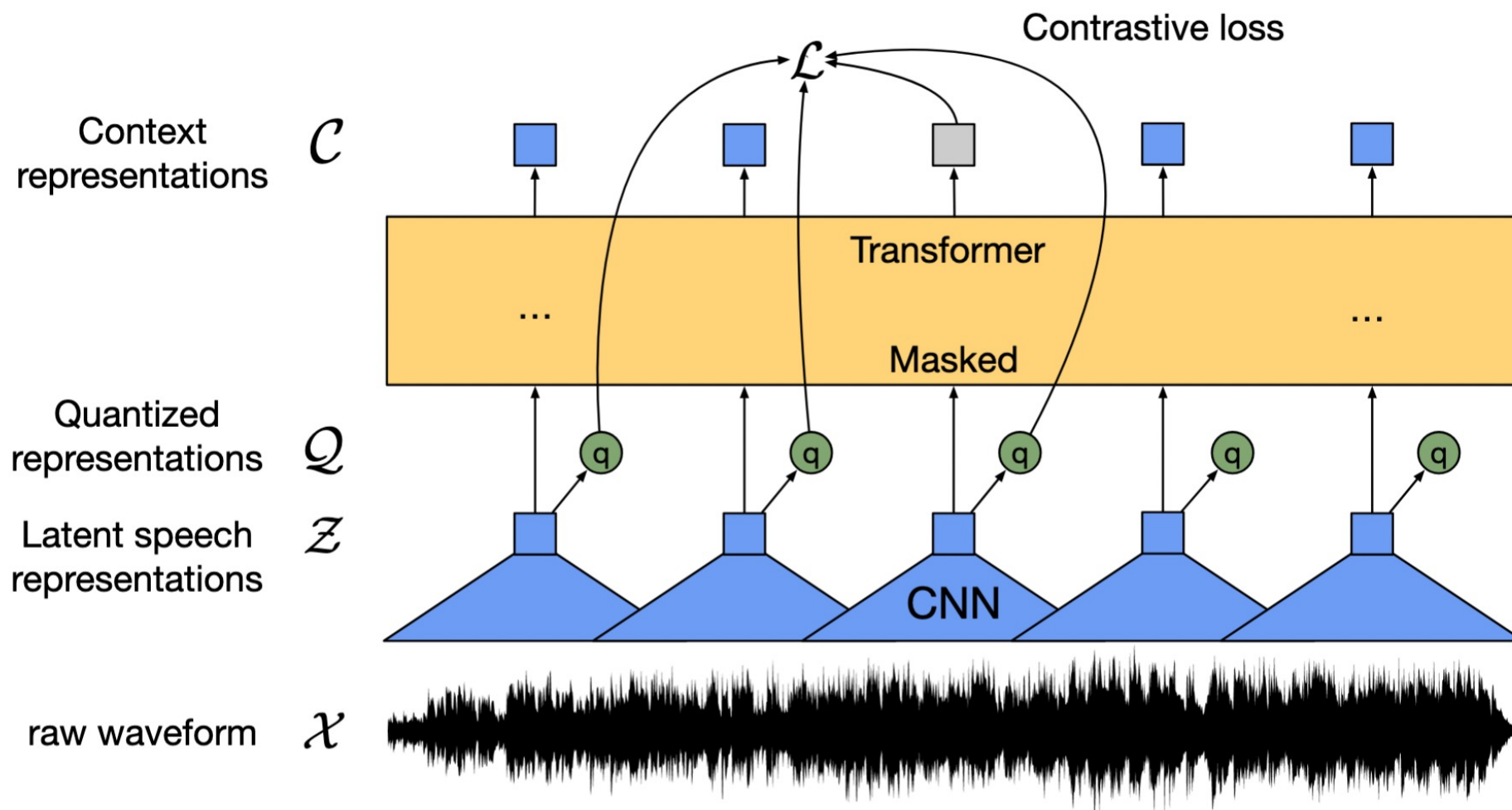
### Latent Feature Encoder (Convolutional Network)

# Hubert (speech resynthesis)

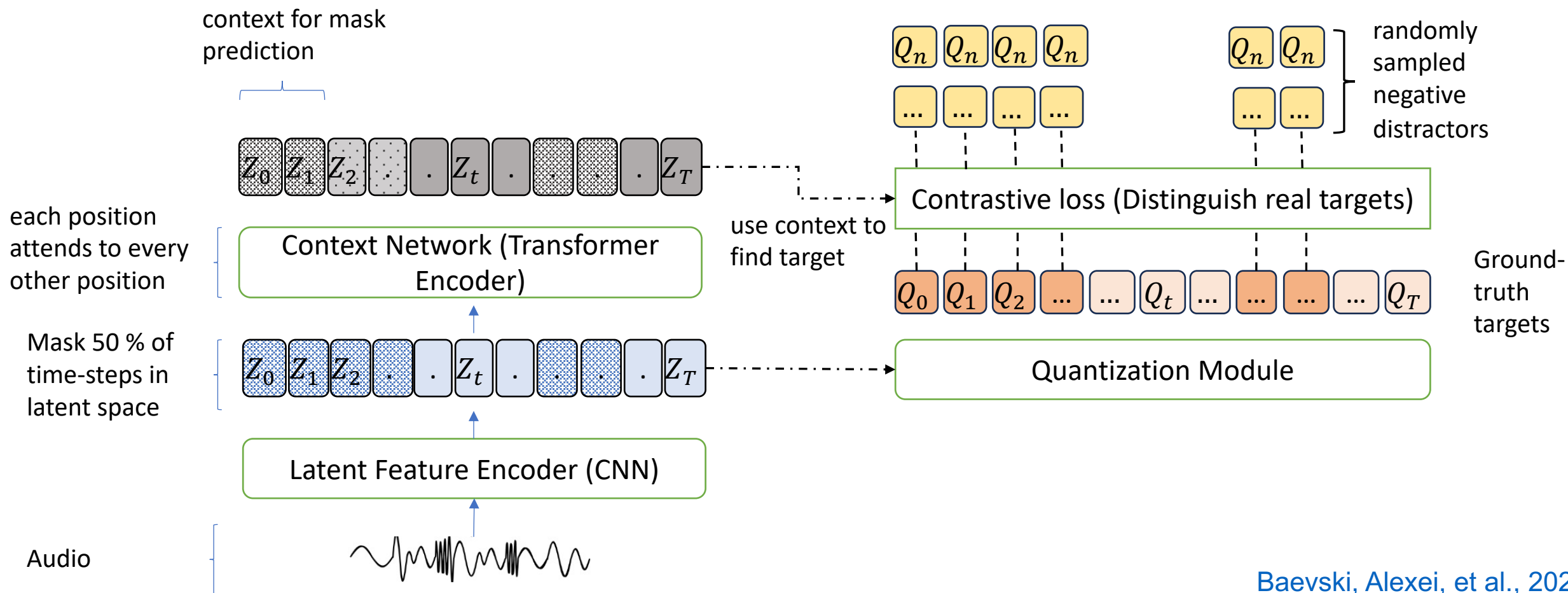




# Wav2Vec 2.0

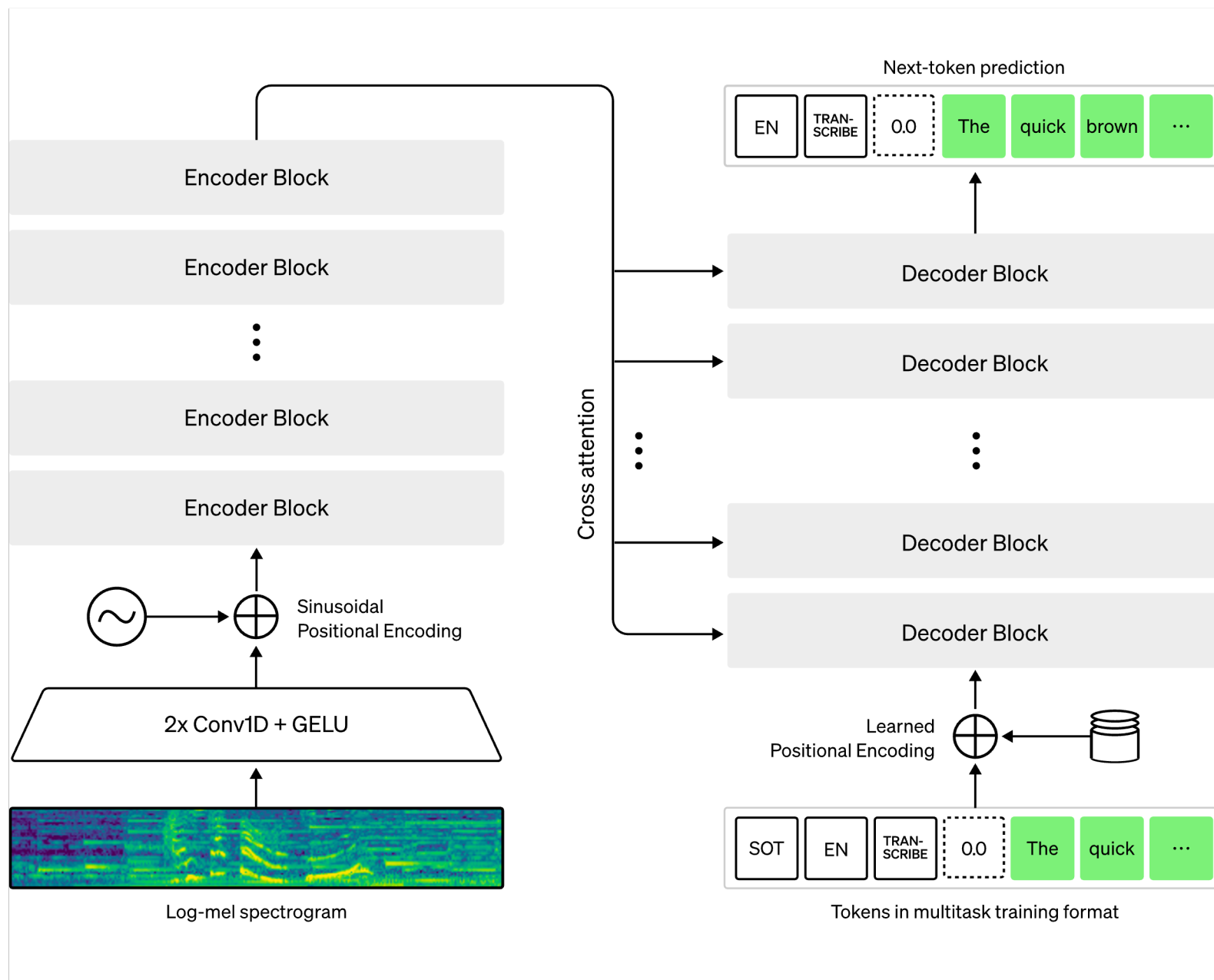


# Wav2Vec 2.0

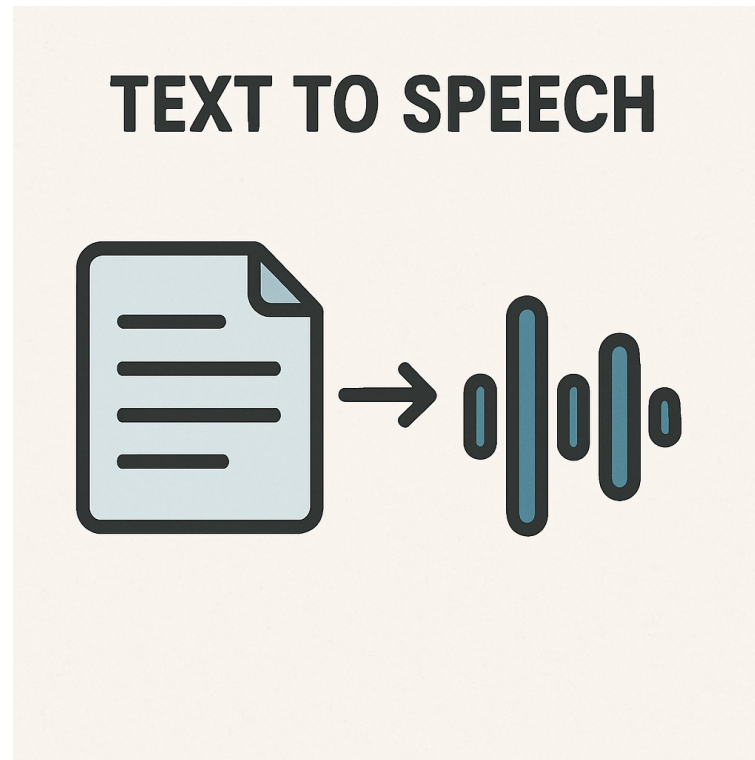


# Whisper

- Trained on 680K hours of multilingual, multitask web data
- Robust to accents, noise, and technical terms
- Supports multilingual transcription and translation to English
- Whisper doesn't outperform models on LibriSpeech but is significantly more robust in zero-shot settings, with 50% fewer errors across varied datasets.



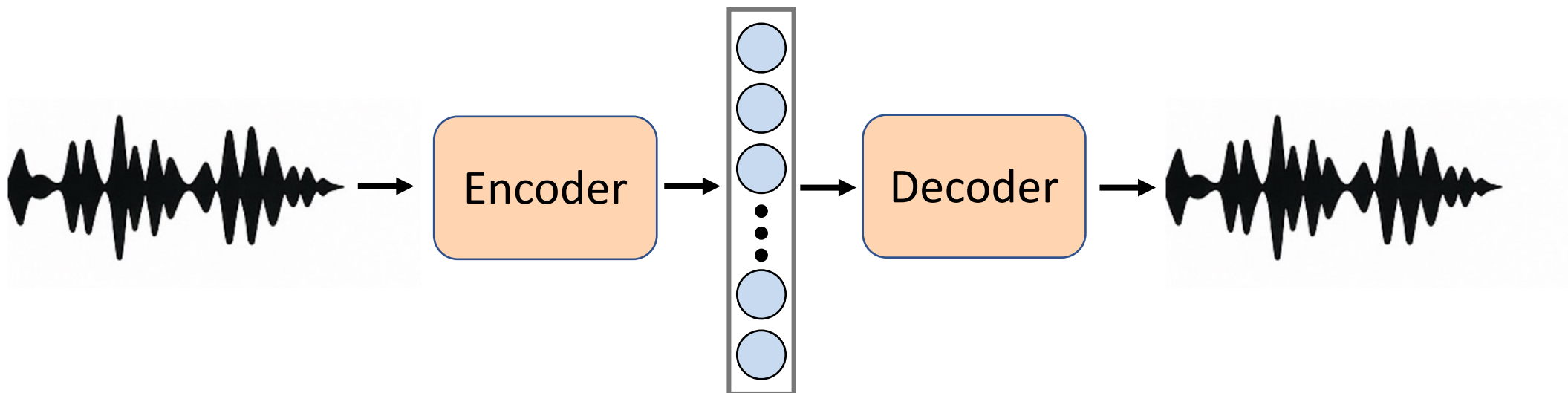
# Thinking multimodal, opens-up possibilities

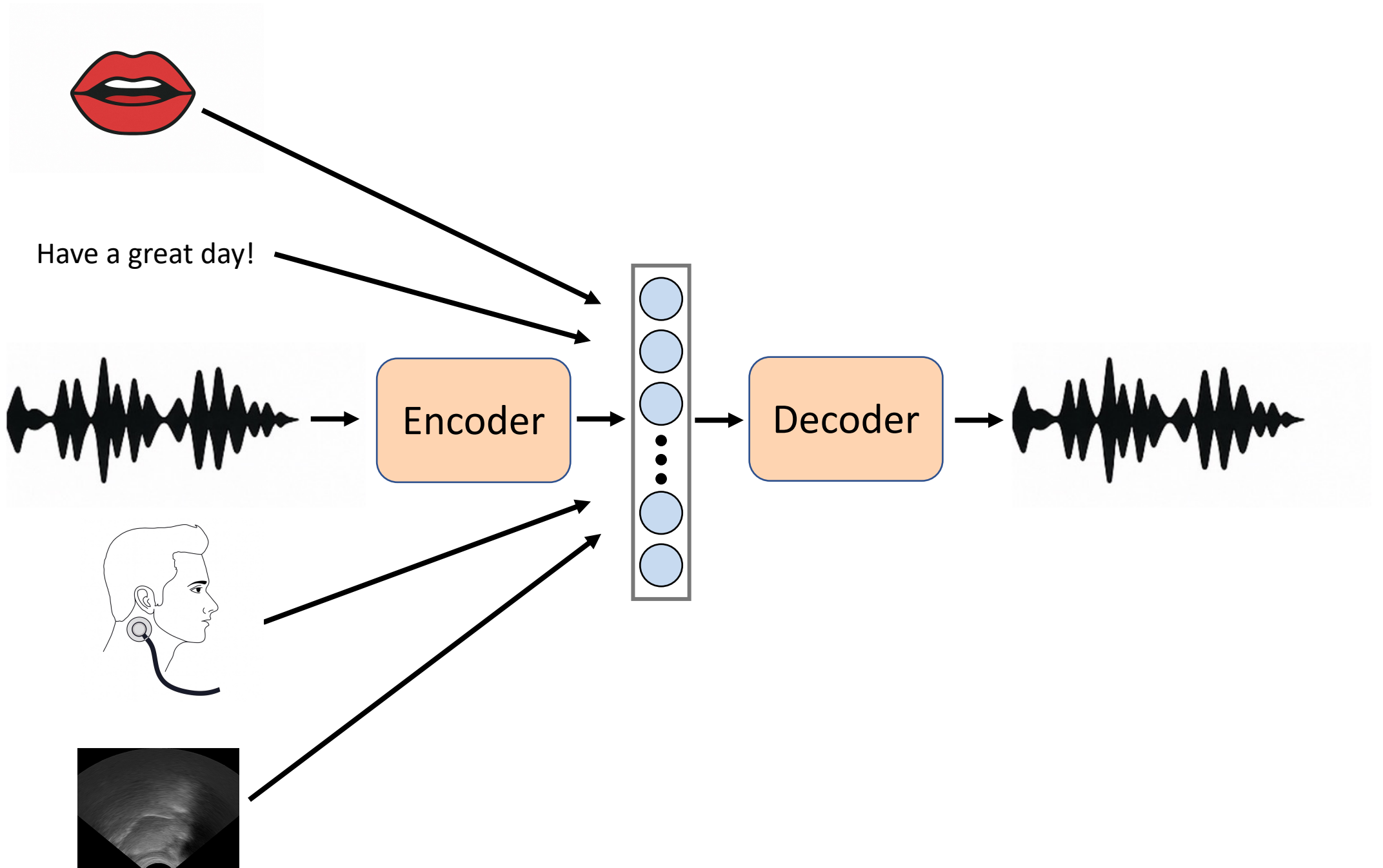




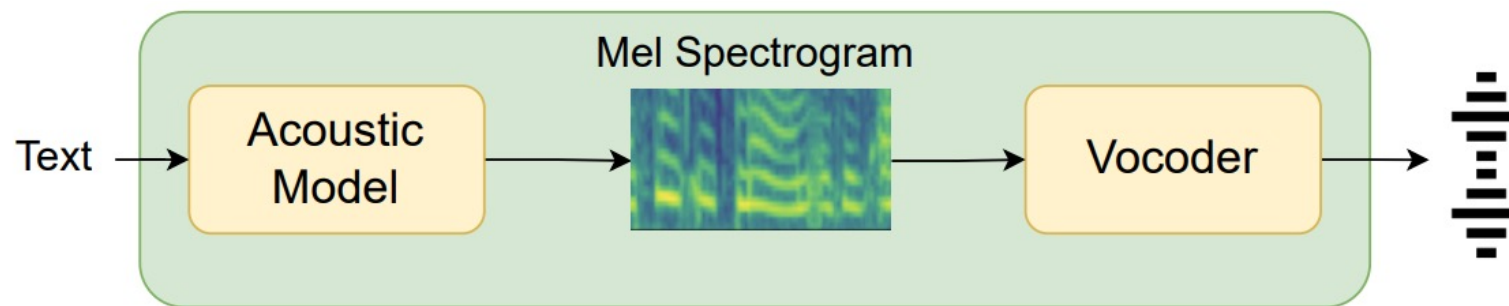
# ParrotTTS

Vocal learning forms the first phase of infants starting to talk (Locke, 1996, 1994) by simply listening to sounds/speech. It is hypothesized (Kuhl and Meltzoff, 1996) that infants listening to ambient language store perceptually derived representations of the speech sounds they hear, which in turn serve as targets for the production of speech utterances. Interestingly, in this phase, the infant has no conception of text or linguistic rules, and speech is considered sufficient to influence speech production (Kuhl and Meltzoff, 1996) as can parrots (Locke, 1994).

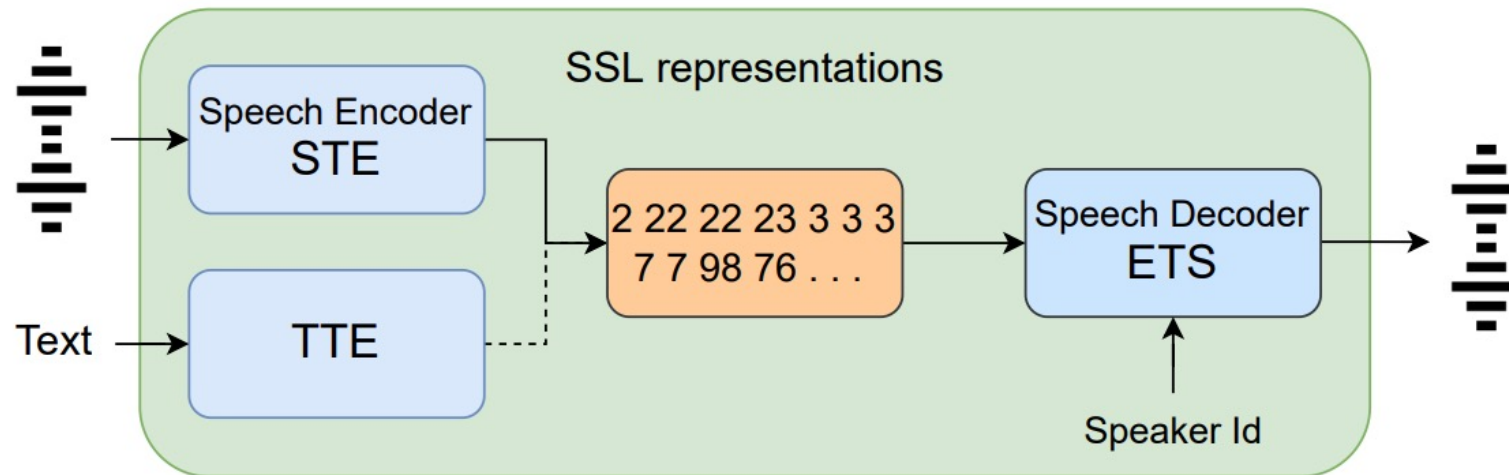




# ParrotTTS



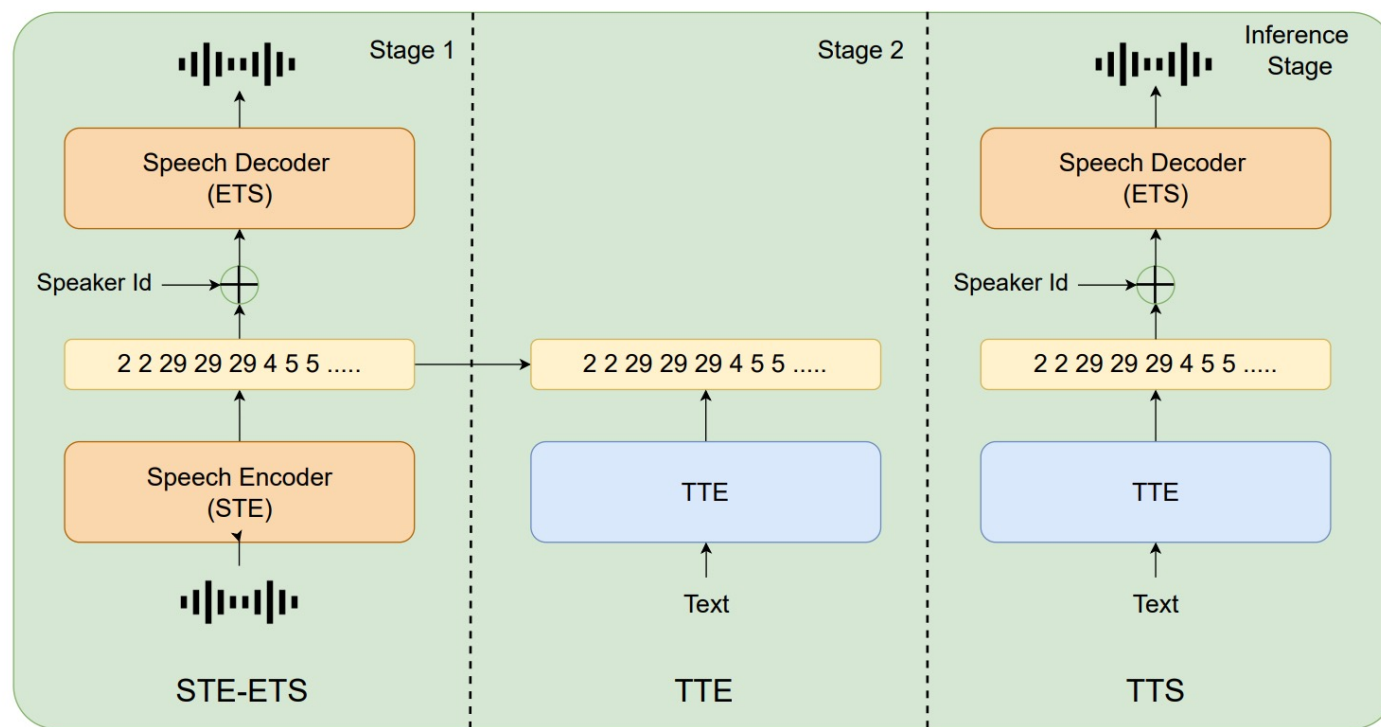
(a)



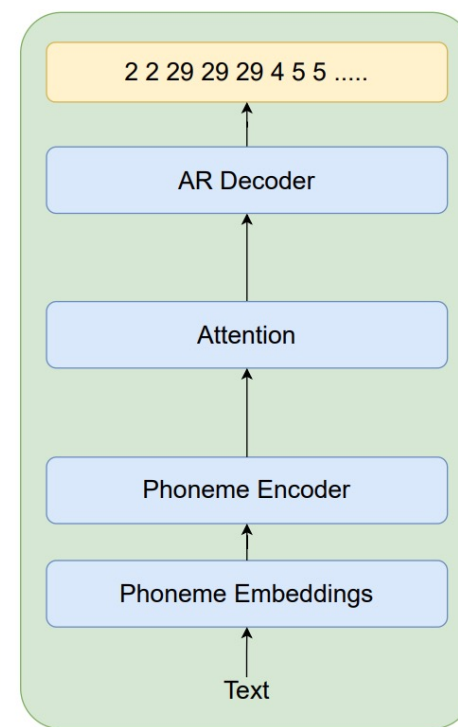
(b)



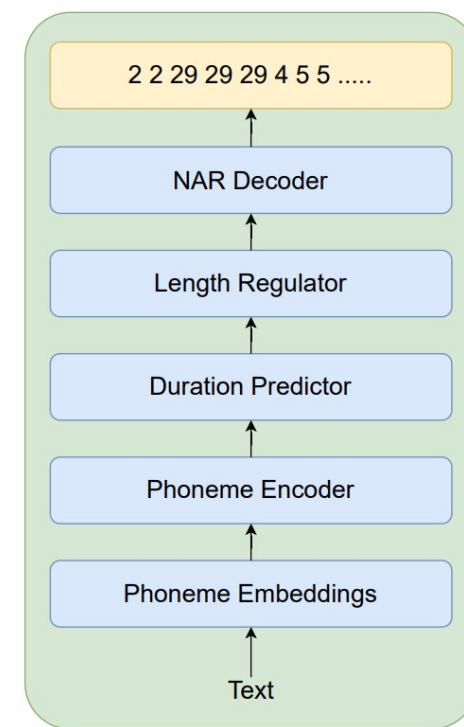
# ParrotTTS



(a) ParrotTTS



(b) Autoregressive TTE



(c) Non Autoregressive TTE

# ParrotTTS

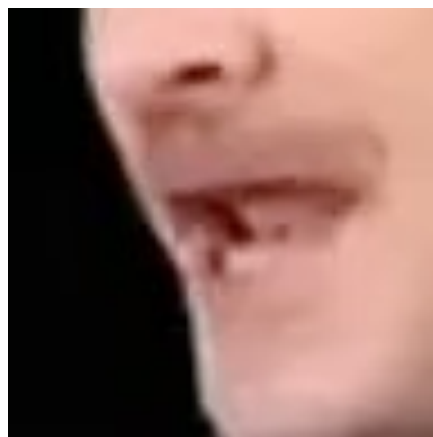
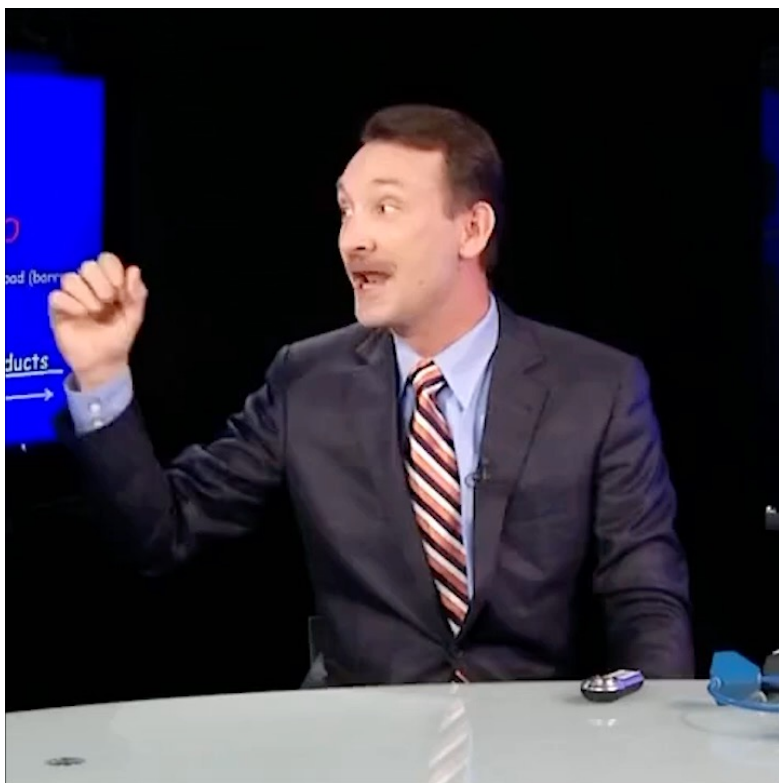
	<b>Model</b>	<b>MOS</b> $\uparrow$	<b>WER</b> $\downarrow$
Traditional TTS	SS-FastSpeech2	3.87	4.52
	SS-Tacotron2	3.90	4.59
	FastSpeech2-SupASR	3.78	4.72
	Tacotron2-UnsupASR	3.50	11.3
WavThruVec	SS-WavThruVec	3.57	6.27
VQ-VAE	SS-VQ-VAES	3.12	21.78
ParrotTTS	AR-TTE <sub>LJS</sub> +SS-ETS	3.85	4.80
	NAR-TTE <sub>LJS</sub> +SS-ETS	3.86	4.58
	NAR-TTE <sub><math>\frac{1}{2}</math>LJS</sub> +SS-ETS	3.81	6.14

Table 1: Subjective and objective comparison of TTS models in the single speaker setting.

<b>Model</b>	<b>VCTK</b>	<b>MOS</b> $\uparrow$	<b>WER</b> $\downarrow$	<b>EER</b> $\downarrow$
GT-Mel+Vocoder	Yes	4.12	2.25	2.12
MS-FastSpeech2	Yes	3.62	5.32	3.21
MS-FastSpeech2-SupASR	No	3.58	6.65	3.85
VC-FastSpeech2	No	3.41	7.44	8.18
WavThruVec-MS	No	3.17	6.79	5.08
NAR-TTE <sub>LJS</sub> +MS-ETS	No	3.78	6.53	4.38

Table 2: Comparison of the multi-speaker TTS models on the VCTK dataset. Column 2 indicates if the corresponding method uses VCTK transcripts while training.

# Lip Reading



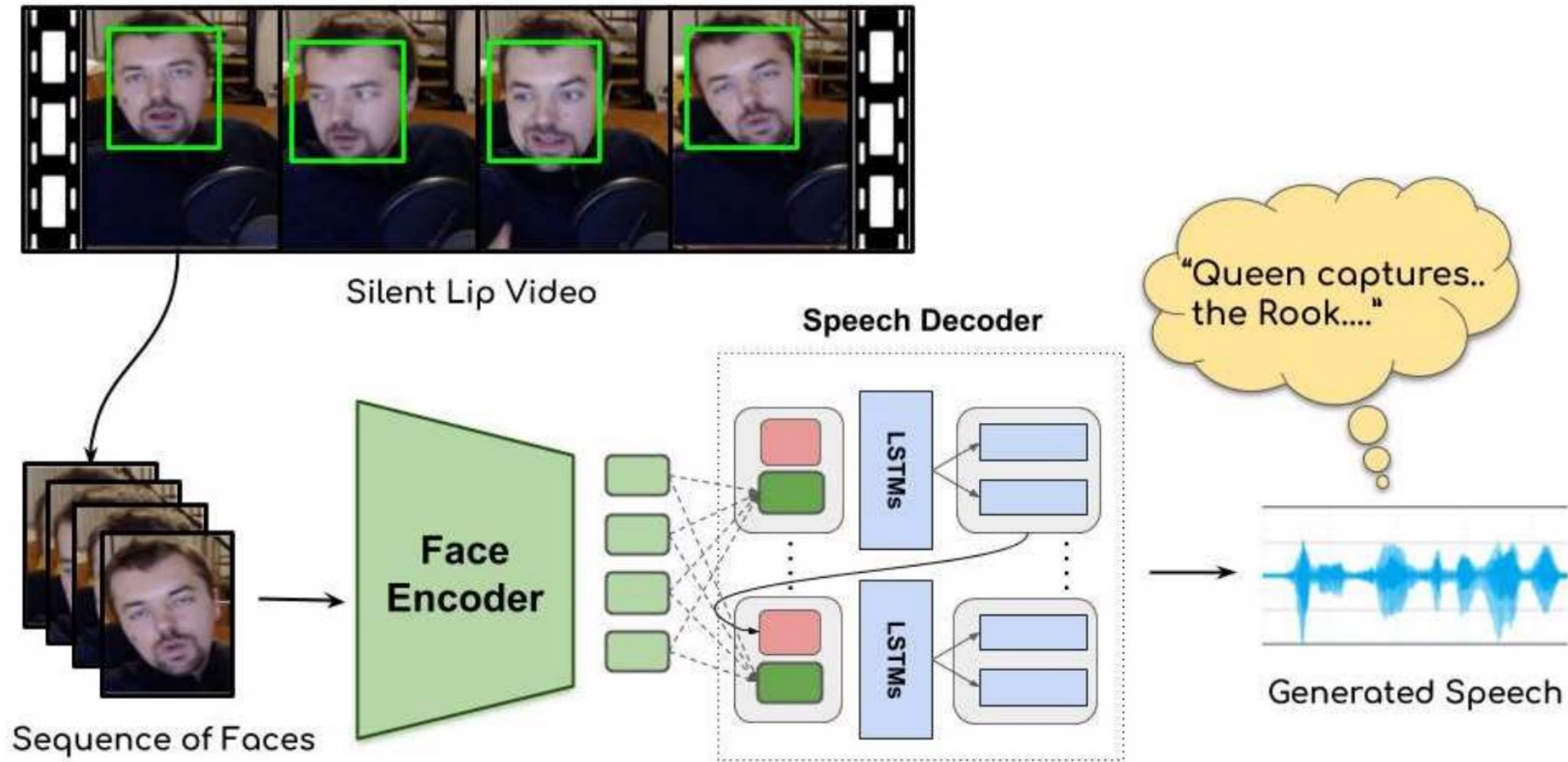
# Lip2Wav

The speech you are hearing is completely generated from the lip movements

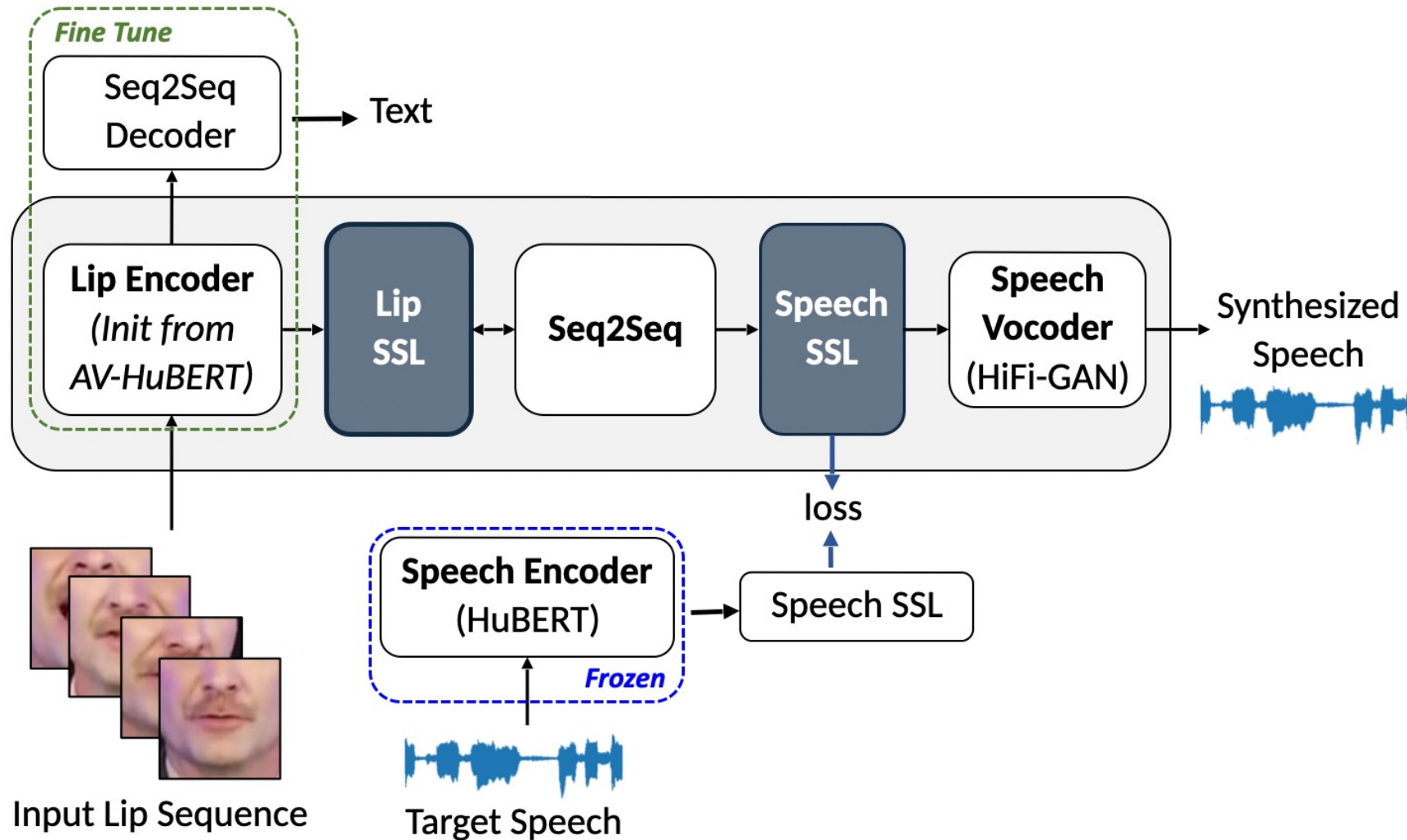




# Lip2Wav



# RobustL2s



# RobustL2S

TABLE II

PERFORMANCE COMPARISON IN CONSTRAINED-SPEAKER SETTING ON GRID-4S DATASET

Method	STOI $\uparrow$	ESTOI $\uparrow$	WER $\downarrow$
Vid2speech [13]	0.491	0.335	44.92 %
Lip2AudSpec [15]	0.513	0.352	32.51 %
1D GAN-based [17]	0.564	0.361	26.64 %
Vocoder-based [40]	0.648	0.455	23.33 %
Ephrat <i>et al.</i> [14]	0.659	0.376	27.83 %
Lip2Wav [5]	0.731	0.535	14.08 %
VAE-based [16]	0.724	0.540	-
VCA-GAN [19]	0.724	<b>0.609</b>	12.25 %
kim <i>et al.</i> [28], [46]	0.738	0.579	-
<b>RobustL2S</b>	<b>0.754</b>	0.571	<b>11.21 %</b>

TABLE III

PERFORMANCE COMPARISON IN CONSTRAINED-SPEAKER SETTING ON TCD-TIMIT-3S DATASET

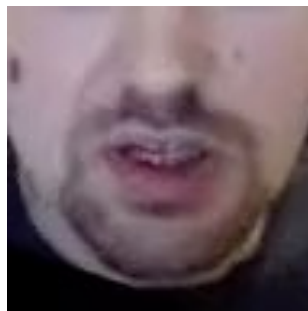
Method	STOI $\uparrow$	ESTOI $\uparrow$	WER $\downarrow$
Vid2speech [13]	0.451	0.298	75.52 %
Lip2AudSpec [15]	0.450	0.316	61.86 %
1D GAN-based [17]	0.511	0.321	49.13 %
Ephrat <i>et al.</i> [14]	0.487	0.310	53.52 %
Lip2Wav [5]	0.558	0.365	31.26 %
VCA-GAN [19]	0.584	0.401	-
<b>RobustL2S</b>	<b>0.596</b>	<b>0.452</b>	<b>29.03 %</b>

TABLE IV

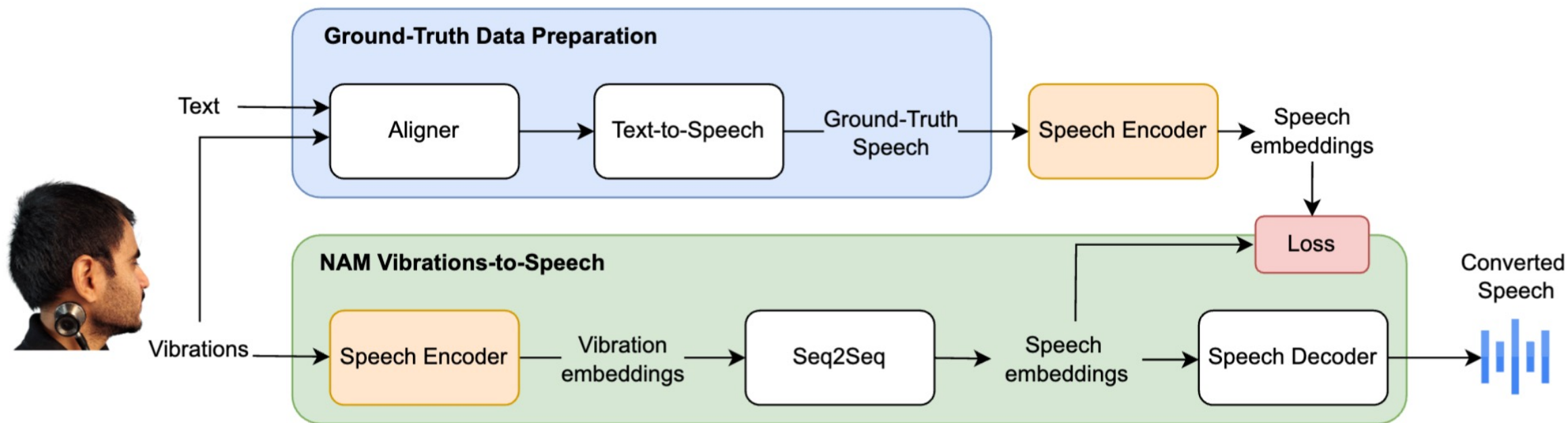
PERFORMANCE COMPARISON IN SPEAKER-DEPENDENT SETTING ON LIP2WAV DATASET

Speaker	Method	STOI $\uparrow$	ESTOI $\uparrow$
Chemistry Lectures (chem)	Ephrat <i>et al.</i> [5]	0.165	0.087
	GAN-based [47]	0.192	0.132
	Lip2Wav [5]	0.416	0.284
	Hong <i>et al.</i> [28]	0.566	<b>0.429</b>
	<b>RobustL2S</b>	<b>0.583</b>	0.397
Chess Analysis (chess)	Ephrat <i>et al.</i> [5]	0.184	0.098
	GAN-based [47]	0.195	0.104
	Lip2Wav [5]	0.418	0.290
	Hong <i>et al.</i> [28]	0.506	0.334
	<b>RobustL2S</b>	<b>0.517</b>	<b>0.340</b>
Deep Learning (dl)	Ephrat <i>et al.</i> [5]	0.112	0.043
	GAN-based [47]	0.144	0.070
	Lip2Wav [5]	0.282	0.183
	Hong <i>et al.</i> [28]	0.576	0.402
	<b>RobustL2S</b>	<b>0.627</b>	<b>0.419</b>
Hardware Security (hs)	Ephrat <i>et al.</i> [5]	0.192	0.064
	GAN-based [47]	0.251	0.110
	Lip2Wav [5]	0.446	0.311
	Hong <i>et al.</i> [28]	0.504	0.337
	<b>RobustL2S</b>	<b>0.511</b>	<b>0.337</b>
Ethical Hacking (eh)	Ephrat <i>et al.</i> [5]	0.143	0.064
	GAN-based [47]	0.171	0.089
	Lip2Wav [5]	0.369	0.220
	Hong <i>et al.</i> [28]	0.463	<b>0.304</b>
	<b>RobustL2S</b>	<b>0.493</b>	0.277

# RobustL2S



# StethoSpeech



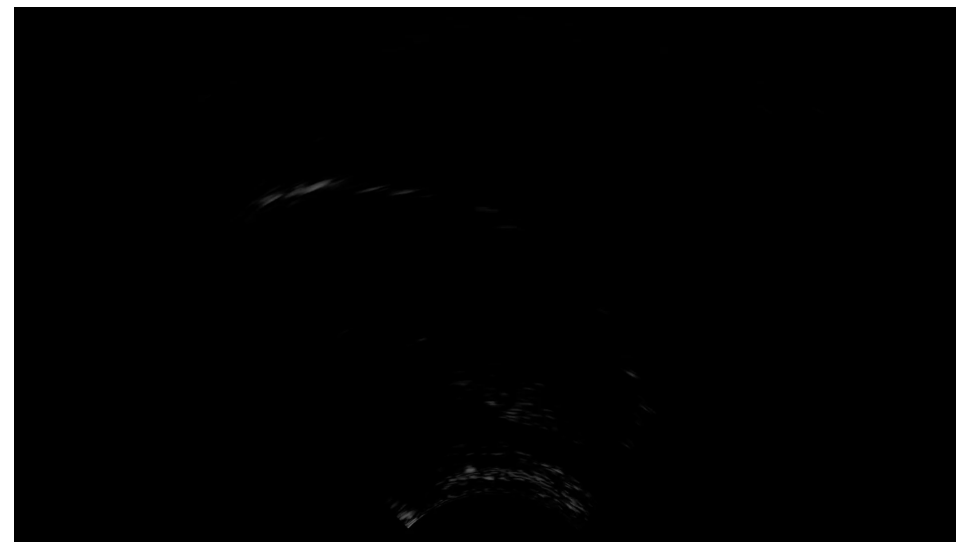


# StethoSpeech

**StethoSpeech: Speech generation through a clinical stethoscope attached to the skin**



# Tongue Ultrasound to speech



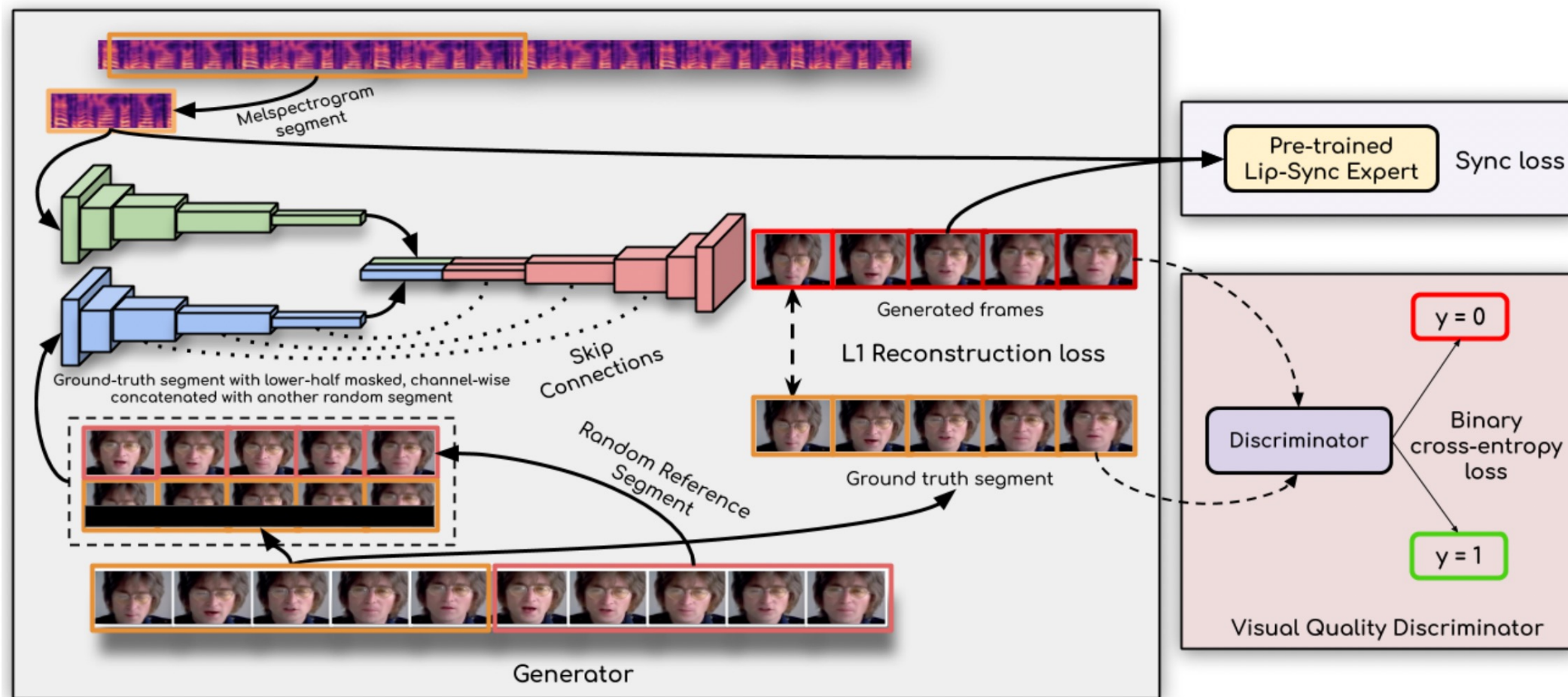
Prediction



Ground Truth

Many fun/useful ideas and possibilities

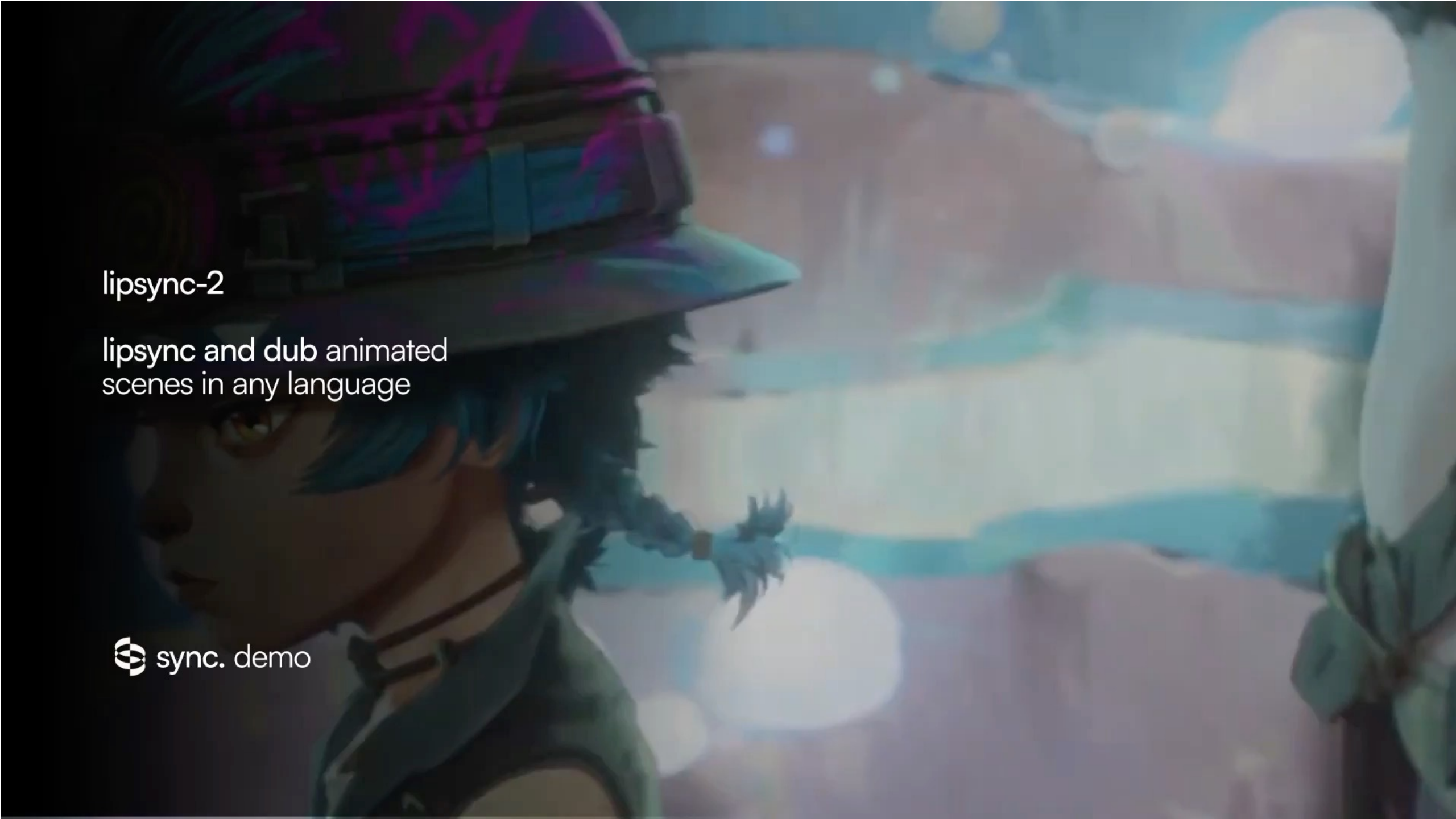
# Wav2Lip



preserve speaker style  
while lipsyncing








lipsync-2

lipsync and dub animated  
scenes in any language

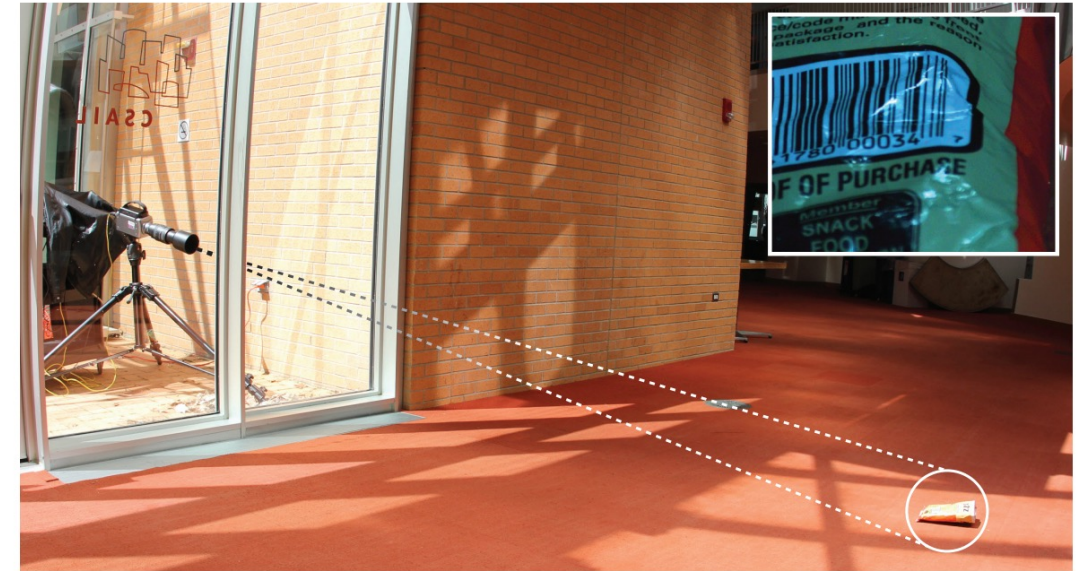
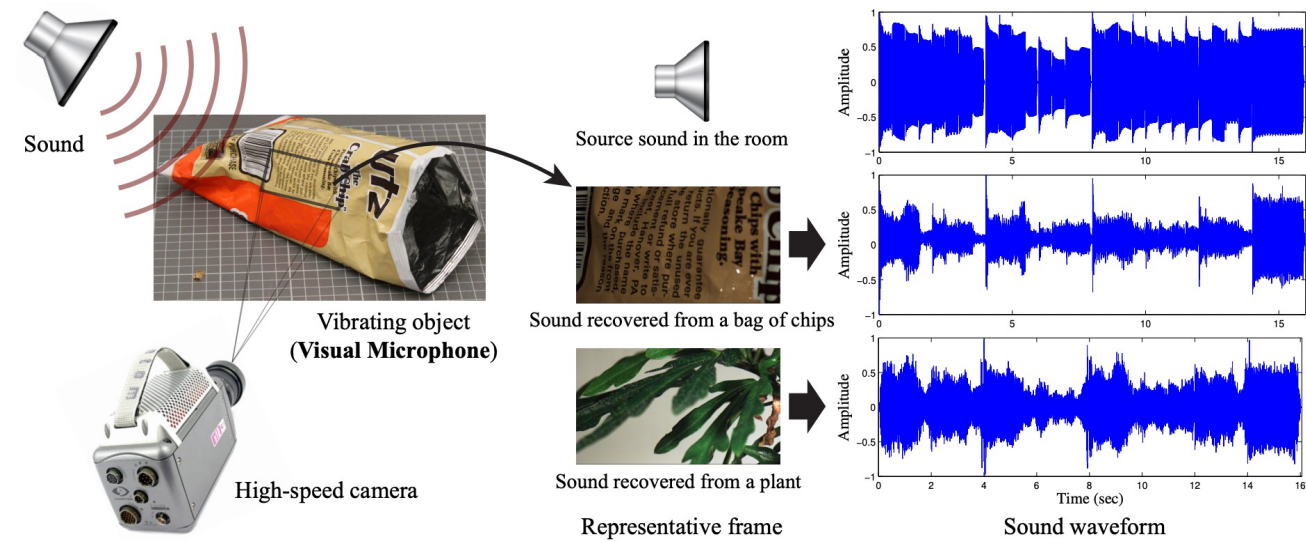
 sync. demo

# Avatars: Synthesia and others

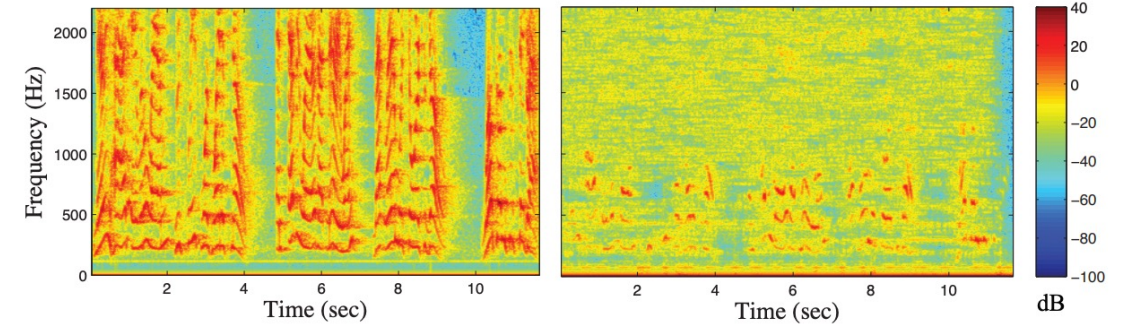




# Visual microphone



(a) Setup and representative frame



(b) Input sound

(c) Recovered sound

# **The Visual Microphone: Passive Recovery of Sound from Video**

**Abe Davis  
Michael Rubinstein  
Neal Wadhwa  
Gautham J. Mysore  
Fredo Durand  
William T. Freeman**

# Automated Cinematography



Original (Wide Angle Static)



Edited



# AudioSet

There are 2,084,320 YouTube videos containing 527 labels

Type a sound to filter the dataset

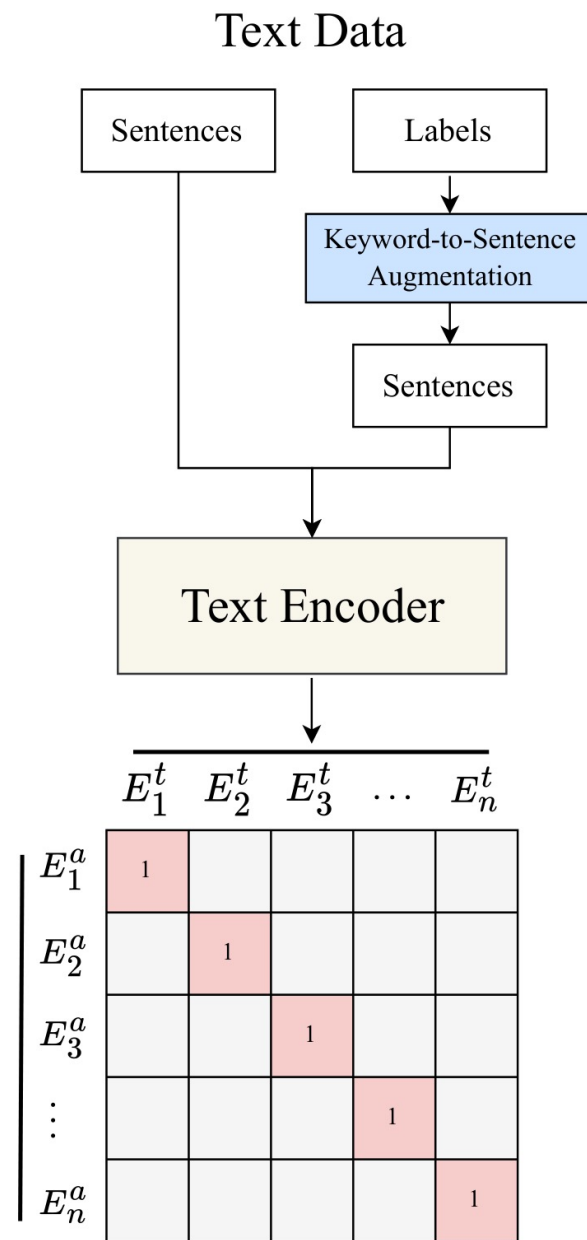
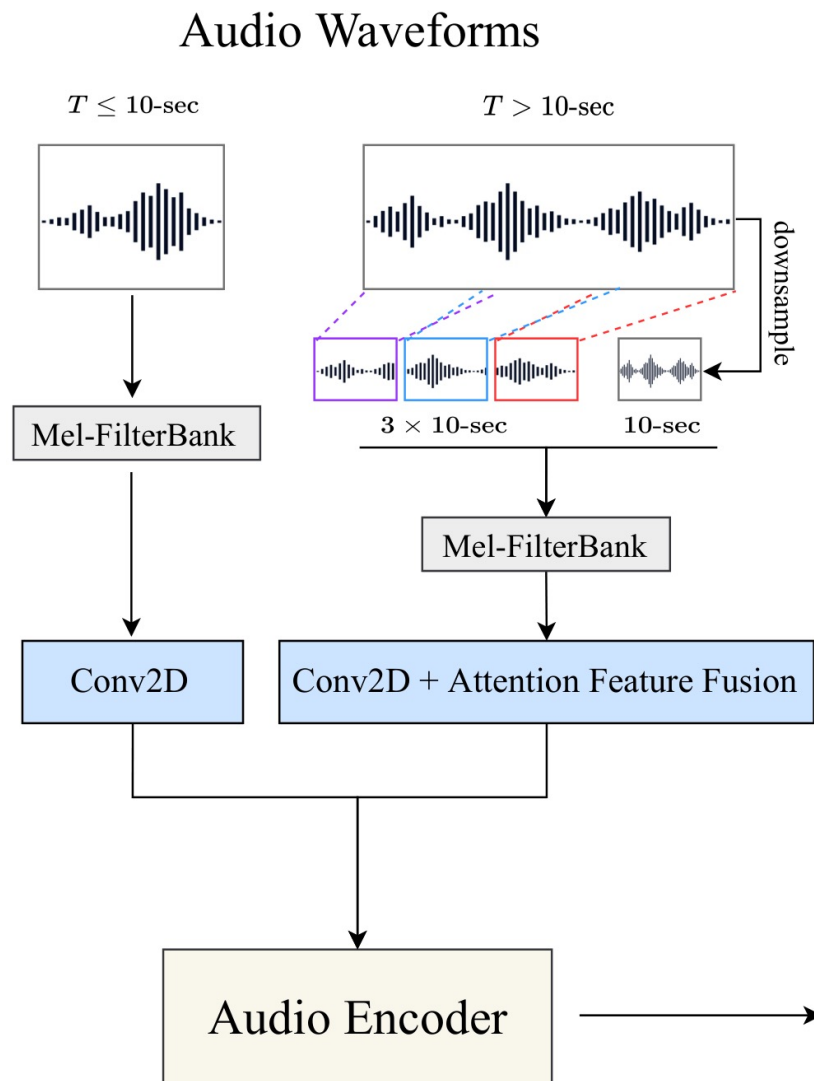
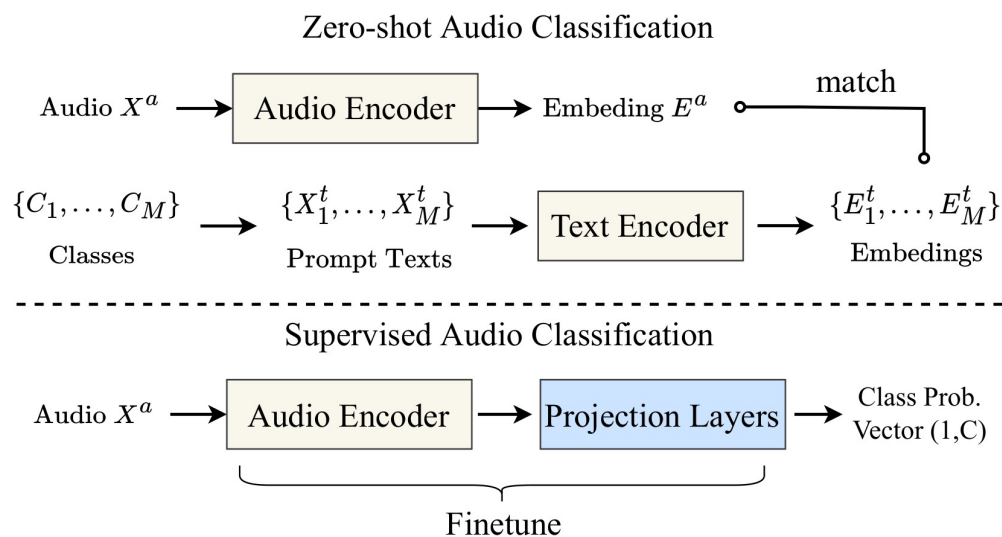


Show detailed breakdown? ☐

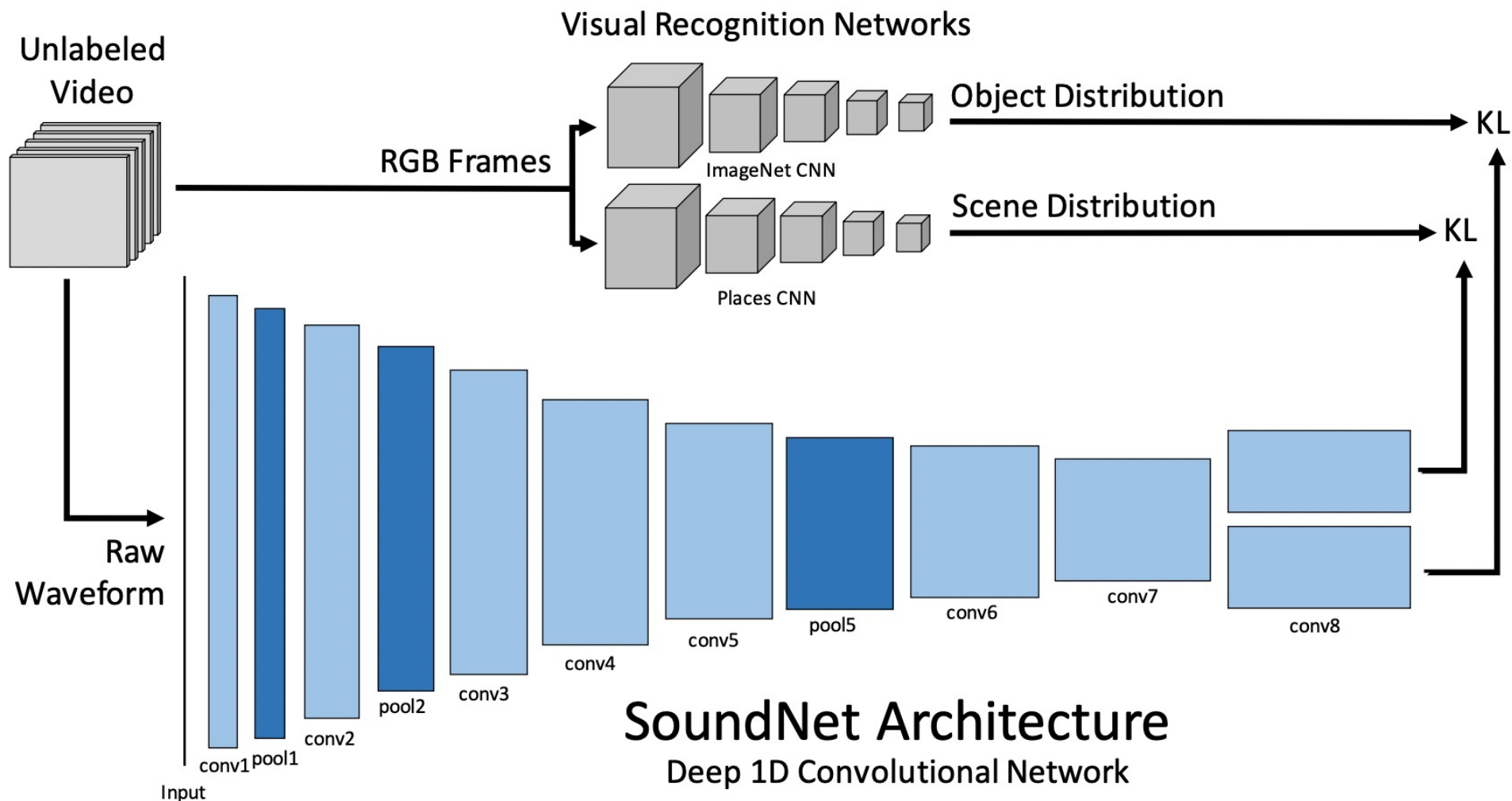
Label	Quality estimate? ▾	Number of videos
<b>Music</b>	100%	1,011,305
<b>Speech</b>	100%	1,010,480
<b>Vehicle</b>	100%	128,051
<b>Musical instrument</b>	100%	117,343
<b>Plucked string instrument</b>	100%	44,565
<b>Singing</b>	100%	42,493
<b>Car</b>	100%	41,554
<b>Animal</b>	100%	40,758
<b>Outside, rural or natural</b>	100%	35,731
<b>Violin, fiddle</b>	100%	28,125
<b>Bird</b>	100%	26,894
<b>Drum</b>	100%	20,246
<b>Engine</b>	100%	16,245
<b>Narration, monologue</b>	100%	15,590
<b>Drum kit</b>	100%	15,169
<b>Acoustic guitar</b>	100%	14,568
<b>Dog</b>	100%	13,705
<b>Child speech, kid speaking</b>	100%	11,816
<b>Bass drum</b>	100%	9,292

<https://research.google.com/audioset/>

# CLAP model



# Multimodal distillation



# Thank you!

