

Conversational AI

MARIA TELEKI

Slides are
posted!

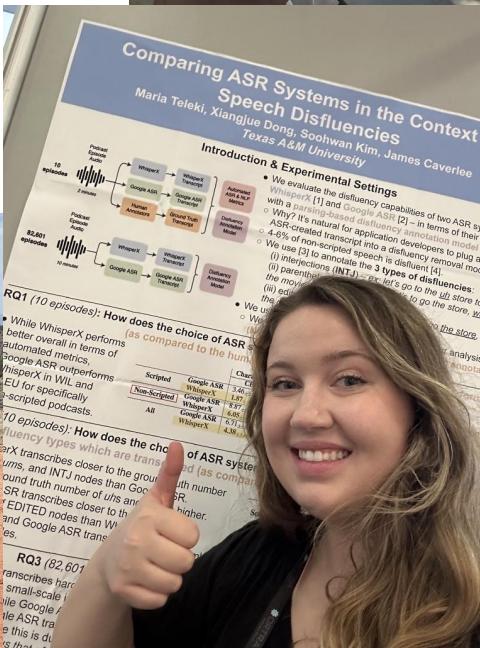


Texas A&M University
Department of Computer Science & Engineering





about me



Agenda

- 1. Why speech?**
- 2. Speech language models**
- 3. Robustness of pipeline approaches**

why speech?

**Speech applications
are everywhere!**



Set a timer for seven, uh, minutes.



Find me the nearest **gas station** — no, sorry, **grocery store**.



Where is my wallet, **you know**, the black one I always carry?

Smart Speaker

Smart Phone

Smart Glasses

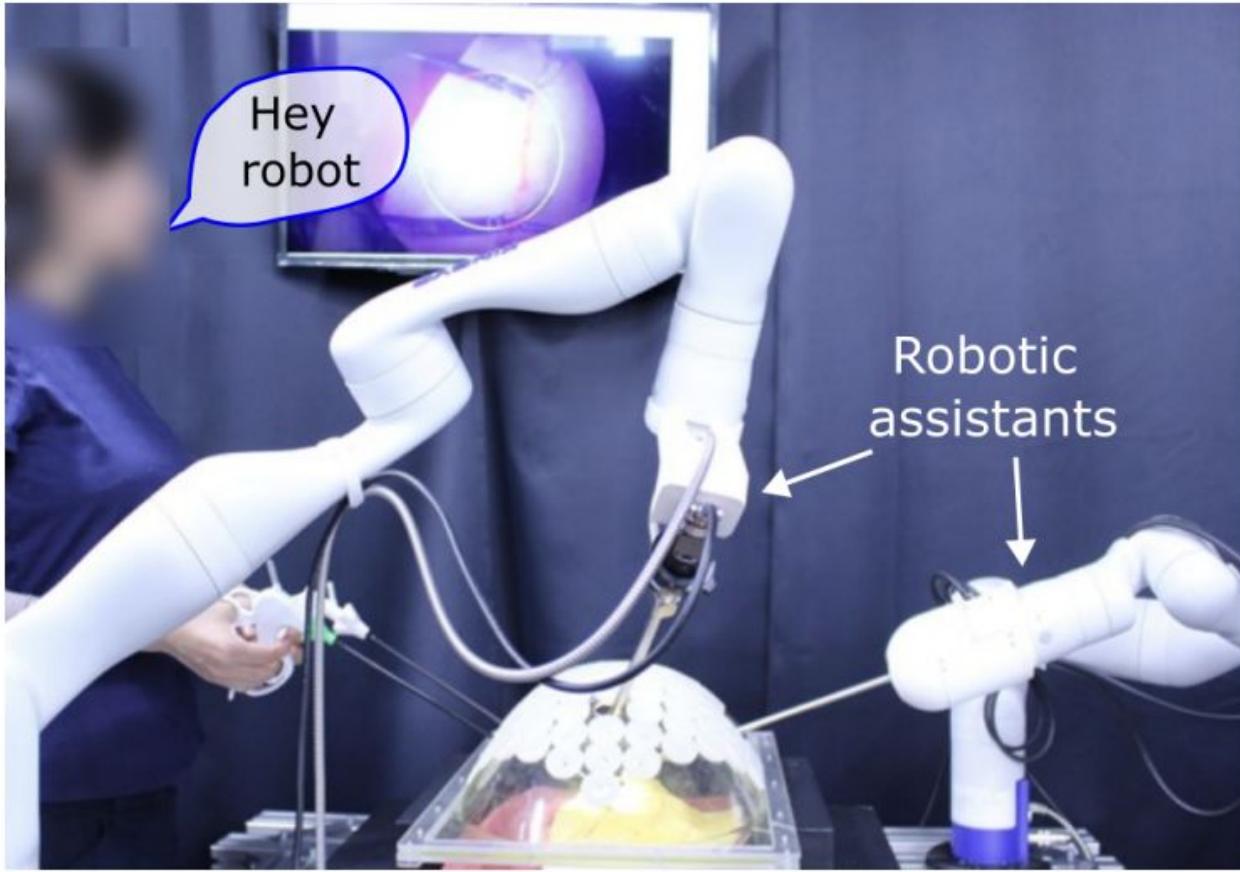


Fig. 1: Voice commands can be used for commanding robotic surgical assistants.

Development Considerations for Implementing a Voice-Controlled Spacecraft System

George Salazar, *NASA Johnson Space Center, Houston, Texas 77058*



■ **Figure 1** CIMON, an early example of an intelligent personal assistant deployed aboard the International Space Station. It served as a proof of concept for voice interaction in microgravity but was limited by its reliance on narrow AI and Earth-based processing.

Advancing Intelligent Personal Assistants for Human Spaceflight:
Leonie Bensch (MIT) Media Lab, Massachusetts Institute of Technology, Cambridge, MA, USA, Oliver Bensch Technical University of Munich (TUM), Germany, Tommy Nilsson N European Space Agency (ESA), Cologne, Germany.

Question w/ your neighbors:

What other speech
applications can you think
of?

Conversational Interfaces are usually called...

CUIs =

Conversational User Interfaces

VUIs =

Voice User Interfaces

The silent conversation device





<https://ruidongzhang.com/research/HPSpeech>

One-sentence summary: HPSpeech turns your over-the-ear headphones into a silent speech interface.

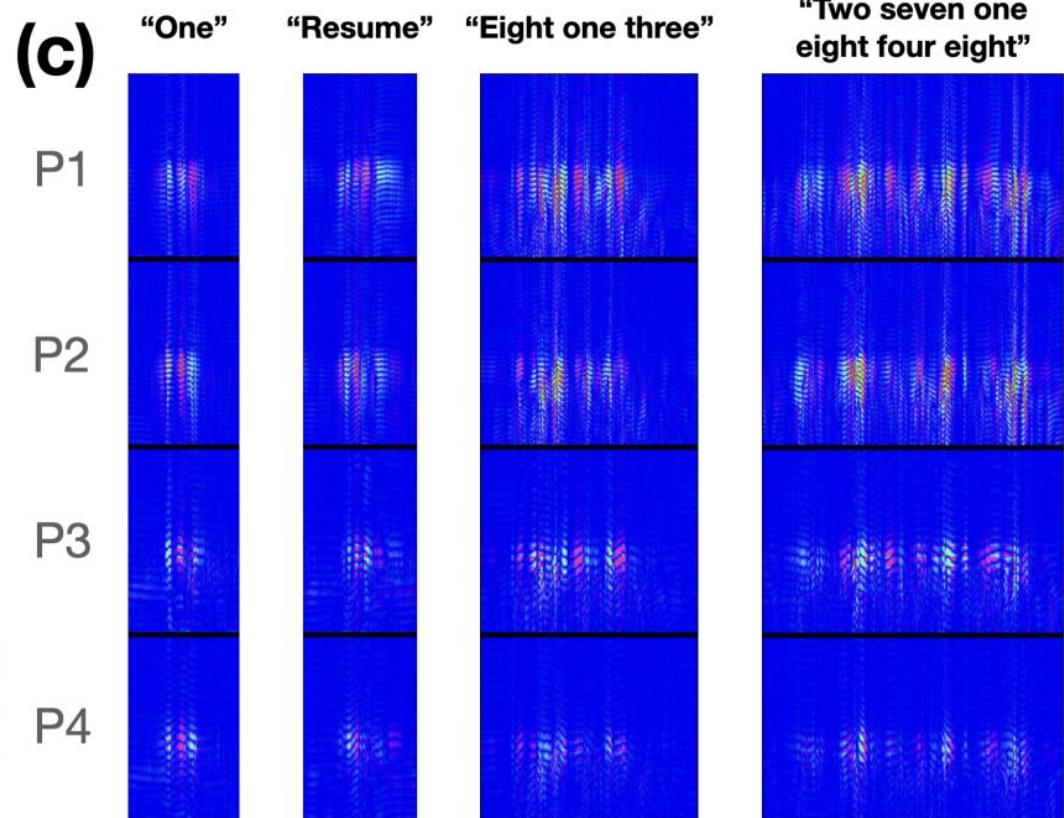
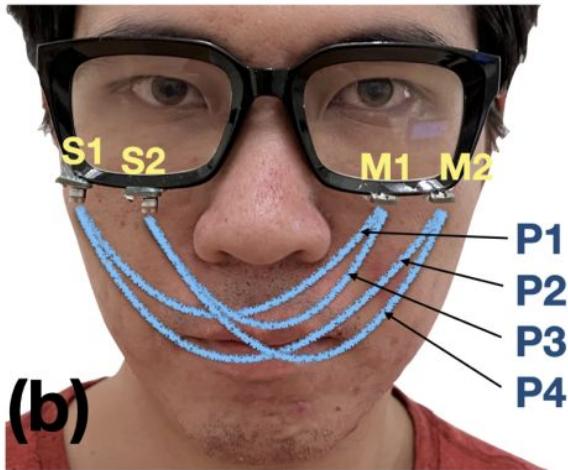
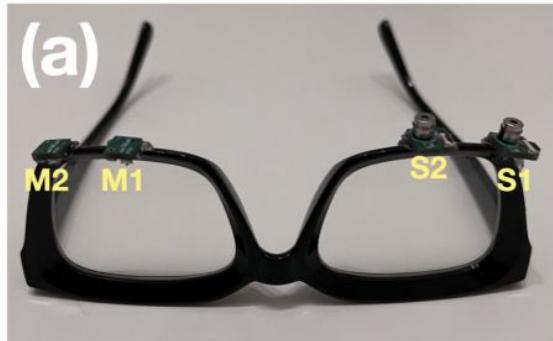


Figure 1: System layout and echo profiles. (a-b)Final sensor position and signal paths. S1, S2: speakers; M1, M2: microphones. P1-P4: Paths. Note that each path consists of multiple path reflection and diffraction that originates from the source speaker and ends at the target microphone. The lines in the figure only illustrate the sources and targets. (c) Echo profiles for different utterances.

Question w/ your neighbors:

Which SSI (*Silent Speech Interface*)
do you think is the best?

Why?

Ok now,

**Let's review a recent
survey of Speech
Language Models.**

The way we speak and write are different.

The way we speak and write are different.

So he calls me up, and he's like, 'I still love you,' and I'm like, I'm just, I mean, this is exhausting, you know – like we are never getting back together. Like, ever.

- TAYLOR SWIFT

The way we speak and write are different.

So he calls me up, and he's like, 'I still love you,' and I'm like, I'm just, I mean, this is exhausting, you know – like we are never getting back together. Like, ever.

- TAYLOR SWIFT

On The Landscape of Spoken Language Models: A Comprehensive Survey

Siddhant Arora^{1*} Kai-Wei Chang^{2*} Chung-Ming Chien^{3*} Yifan Peng^{1*} Haibin Wu^{2*#}
Yossi Adi⁴⁺ Emmanuel Dupoux⁵⁺ Hung-Yi Lee²⁺ Karen Livescu³⁺ Shinji Watanabe¹⁺

¹ Carnegie Mellon University, USA

² National Taiwan University, Taiwan

³ Toyota Technological Institute at Chicago, USA

⁴ Hebrew University of Jerusalem, Israel

⁵ ENS - PSL, EHESS, CNRS, France

ArXiv, highlighted at INTERSPEECH '25 keynote

Big picture context is we live in a **MULTIMODAL** world

LLMs = <i>Large Language Models</i>	SLMs = <i>Speech Language Models</i>	VLMs = <i>Vision Language Models</i>	RLMs = <i>LLMs for Recommendation (I made this name up)</i>	...
$p(\text{text} \text{text})$	$p(\text{text} \text{speech},\text{text})$	$p(\text{text} \text{image},\text{text})$	$p(\text{text} \text{collaborative-signal},\text{text})$	

And right now people are building models for all these different types of modalities

SLMs can do tasks that LLMs can't do

Task	Examples of Instructions
Speech recognition	Recognize the speech and give me the transcription. (Tang et al., 2024) Repeat after me in English. (Grattafiori et al., 2024)
Speech translation	Translate the following sentence into English. (Grattafiori et al., 2024) Recognize the speech, and translate it into English (Chu et al., 2023)
Speaker recognition	How many speakers did you hear in this audio? Who are they? (Tang et al., 2024)
Emotion recognition	Describe the emotion of the speaker. (Tang et al., 2024) Can you identify the emotion? Categorise into: sad, angry, neutral, happy (Das et al., 2024)
Question answering	What happened to this person? (Wang et al., 2023b) Generate a factual answer to preceding question (Das et al., 2024) What medicine is mentioned? Briefly introduce that medicine. (Peng et al., 2024b)

Table 2: Examples of instructions for speech-related tasks used in SLM instruction tuning.

SLMs are trained in a super cool new way!!!!!!

Table 1: Typology of text and spoken LMs. We use a loose notation here, where *speech* and *text* are to be interpreted in context; for example, $p(\text{text}|\text{text})$ in post-trained text LMs corresponds to modeling some desired output text given an input text instruction or prompt. “Post-training” refers to any form of instruction-tuning and/or preference-based optimization of the SLM. Please see the sections below for details and references for the example models.

Type of LM	Training Strategy	Model distribution	Examples
pure text LM	pre-training	$p(\text{text})$	GPT, Llama
pure text LM	post-training	$p(\text{text} \text{text})$	ChatGPT, Llama-Instruct
pure speech LM	pre-training	$p(\text{speech})$	GSLM, AudioLM, TWIST
pure speech LM	post-training	$p(\text{speech} \text{speech})$	Align-SLM
speech+text LM	pre-training	$p(\text{text}, \text{speech})$	SpiRit-LM, Moshi (pre-trained)
speech+text LM	post-training	$p(\text{text}, \text{speech} \text{text}, \text{speech})$	Moshi (post-trained), Mini-Omni
speech-aware text LM	post-training	$p(\text{text} \text{speech}, \text{text})$	SALMONN, Qwen-Audio-Chat

“tokenizing speech” → speech becomes like text for modeling



Pipeline Approach vs. End-to-End Approach



✓ preferred by industry for *controllability* → think custom vocabulary & ease of debugging

✓ currently in research stage but very promising → major barrier imo is downsampling problem (will explain)

The big meta-level question to pay attention to...

- What's the input?
- What's the output?
- How do you glue the parts together?
- Based on the input and output, **what is the SIGNAL being learned (by the model)?**

Let's talk SLM architecture

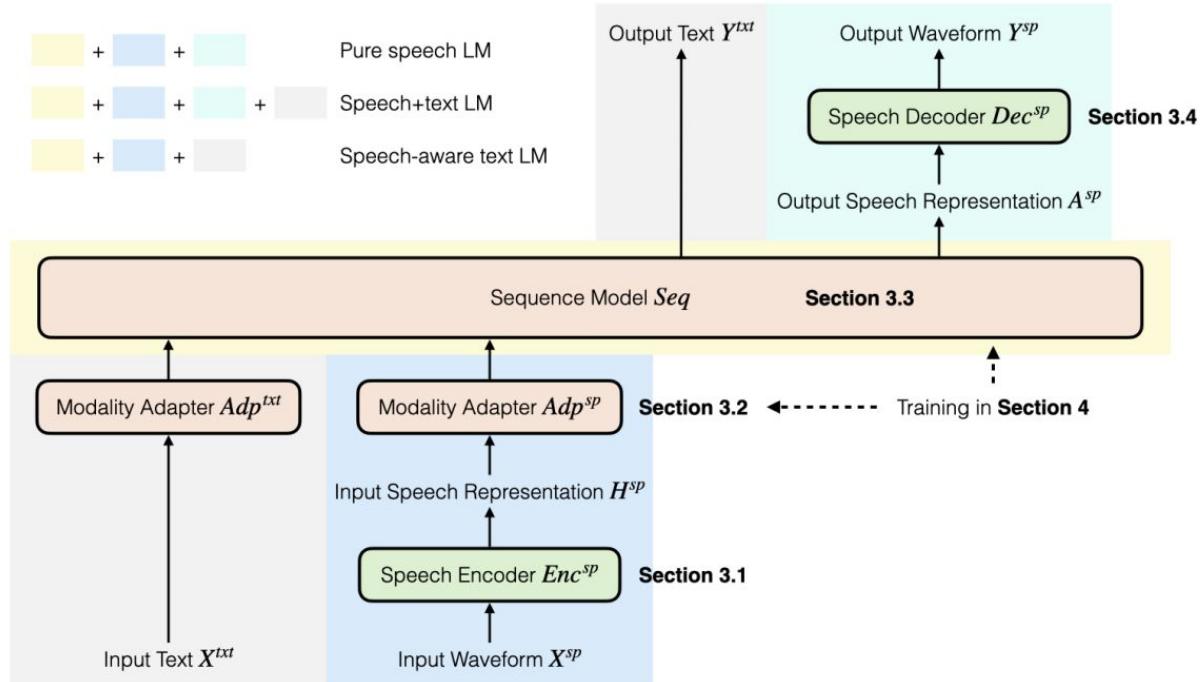


Figure 1: Overview of SLM architecture. See Sections 3 and 4 for more detailed descriptions of the components and training methods, respectively.

The paper is basically walking through this pic like it's a map – so that's what we're going to do! :)

Question w/ your neighbors:

How are SLMs
similar/different from LLMs?

Let's talk SLM architecture

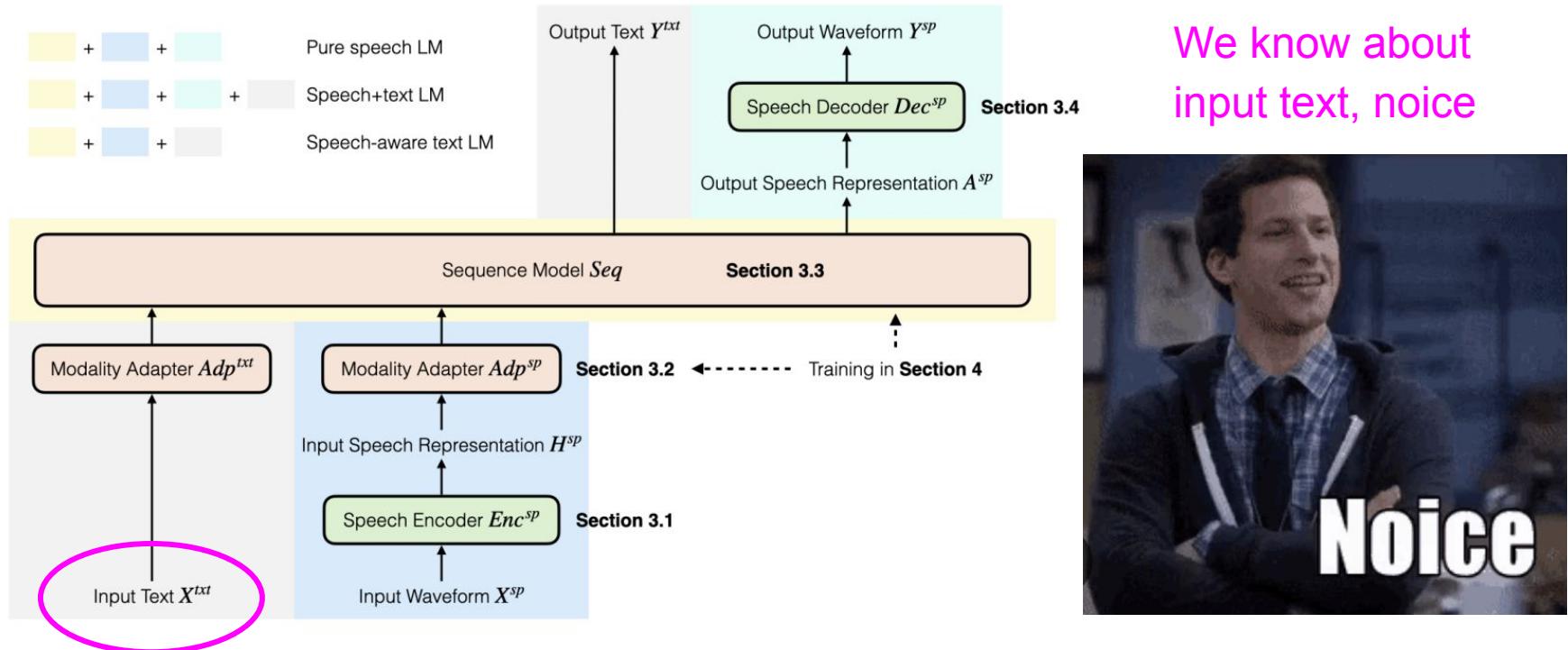
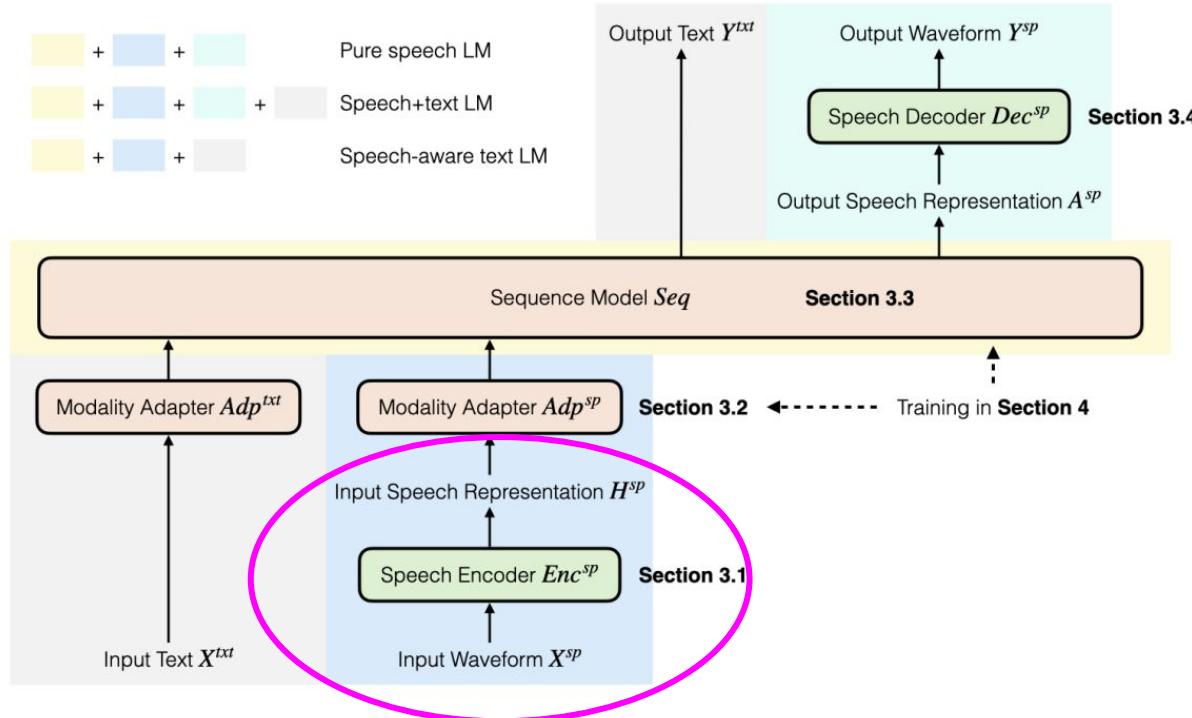


Figure 1: Overview of SLM architecture. See Sections 3 and 4 for more detailed descriptions of the components and training methods, respectively.

Let's talk SLM architecture



So now let's walk through the speech encoder!

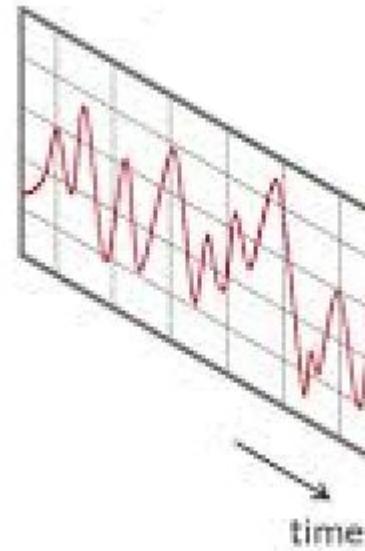
Figure 1: Overview of SLM architecture. See Sections 3 and 4 for more detailed descriptions of the components and training methods, respectively.

Let's start from the speech signal

What is it? It's air pressure over time.

Microphones pick up differences in air pressure – that's sound.

So we get a lil plot of that.



Let's start from the speech signal

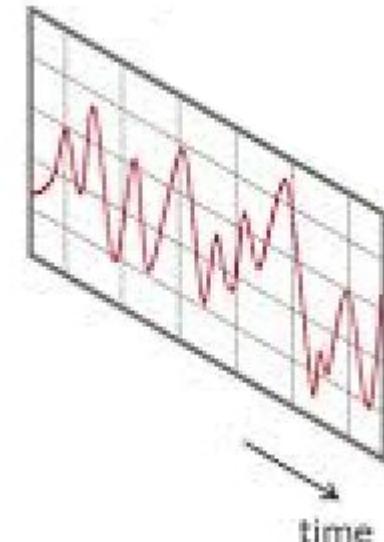
What is it? It's air pressure over time.

Microphones pick up differences in air pressure – that's sound.

So we get a lil plot of that.

What's the problem? Why can't we just use that?

You have to super-duper oversample:



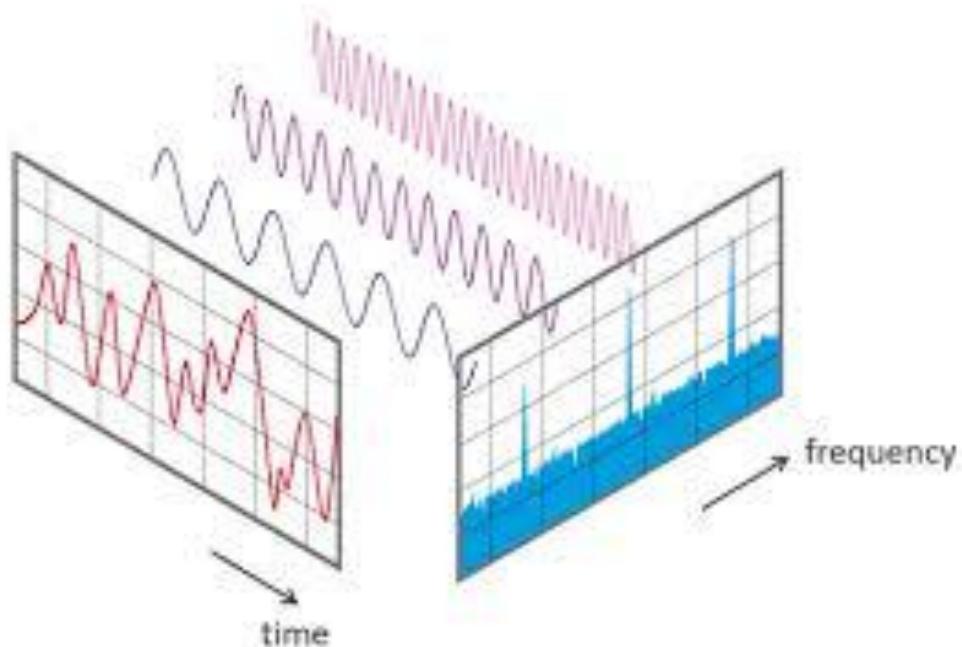
. Redundancy in Raw Waveform

- Oversampling relative to perception:
At 16 kHz, you get 16,000 samples per second. But most speech information lives below ~4 kHz (Nyquist), and much of it is even lower (<2 kHz). → lots of extra samples.
- Correlation across samples:
Neighboring values in the waveform are highly correlated (they don't change randomly). For example, one period of a 200 Hz vowel spans ~80 samples — most of those points contain *the same information about the harmonic structure*.
- Stationarity over frames:
Speech characteristics (formants, pitch) are stable for 20–50 ms. That's **hundreds of samples** where the "meaning" doesn't change.
 - o raw data is very long, repetitive, and inefficient to model directly.

But, we have a trick!

We can do a fourier transform –

Just pull out the frequency components for that time chunk!

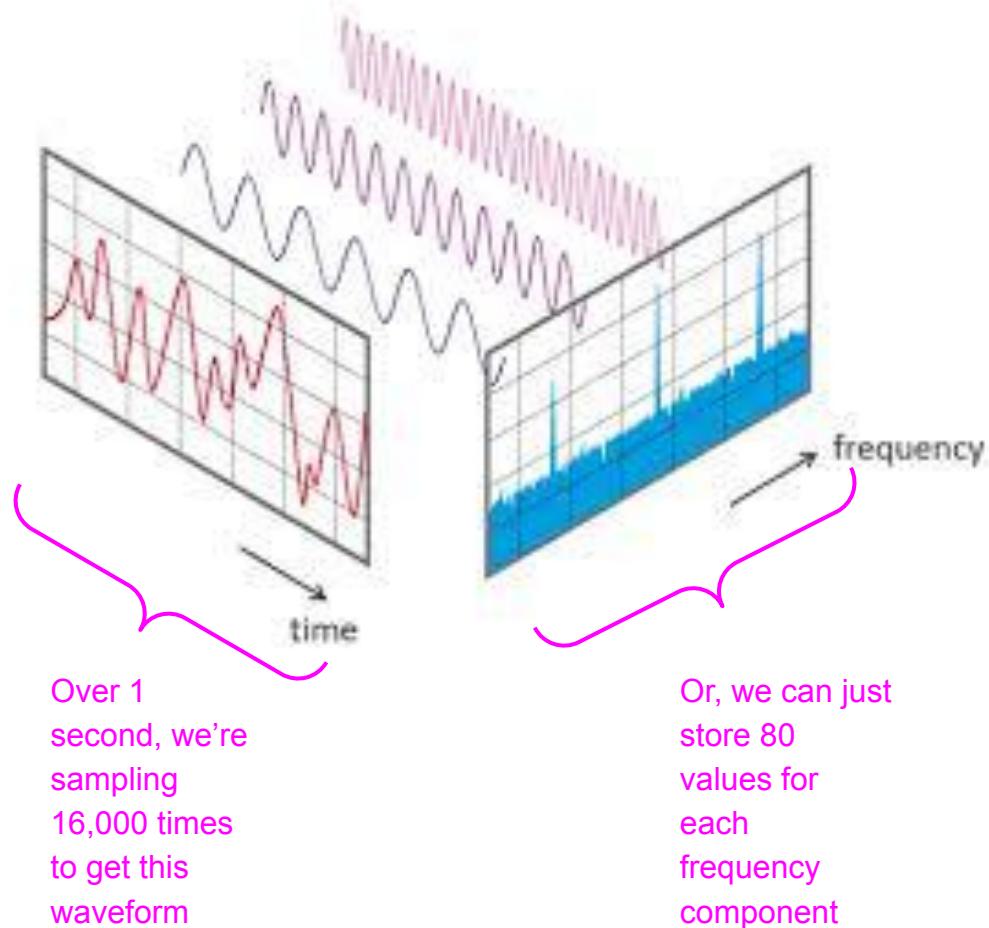


But, we have a trick!

We can do a fourier transform –

Just pull out the frequency components for that time chunk!

So now we can store 80 values per second *instead of 16,000 values per second*



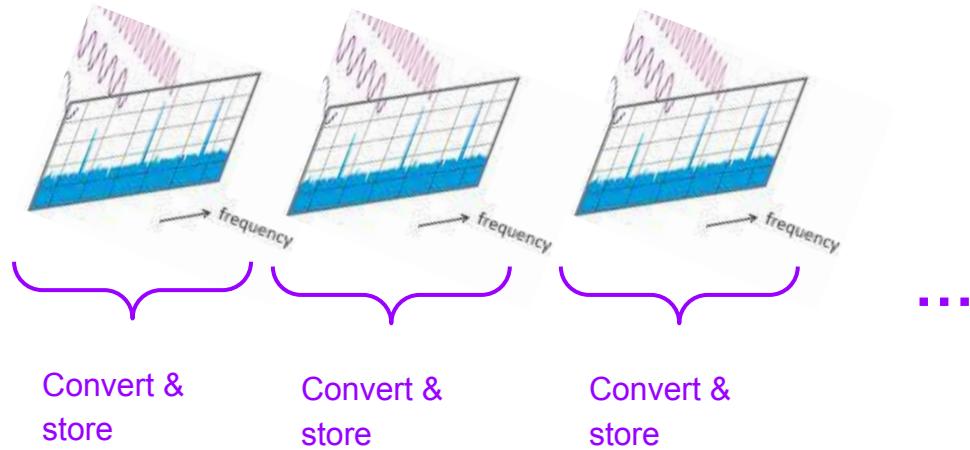
But, we have a trick!

We can do a fourier transform –

Just pull out the frequency components for that time chunk!

So now we can store 80 values per second *instead of 16,000 values per second*

So we store a “frequency snapshot” each second instead



But, we have a trick!

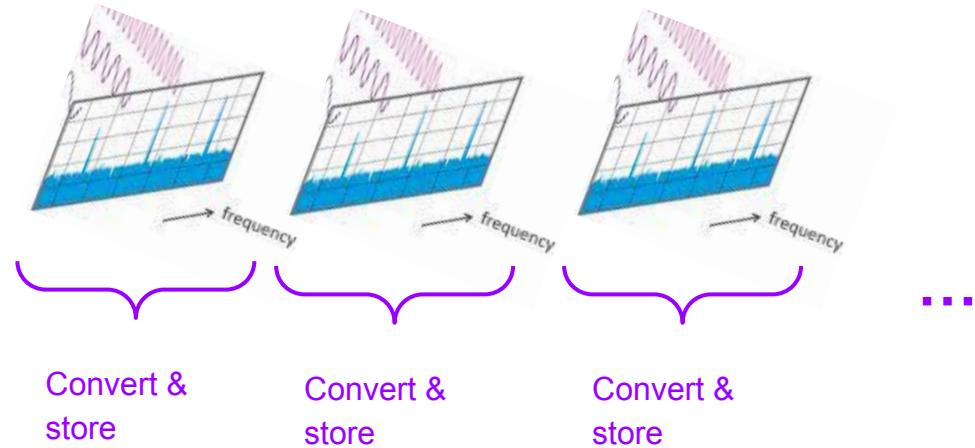
We can do a fourier transform –

Just pull out the frequency components for that time chunk!

So now we can store 80 values per second *instead of 16,000 values per second*

So we store a “frequency snapshot” each second instead

→ so we track how frequencies change over *time* while massively reducing data.

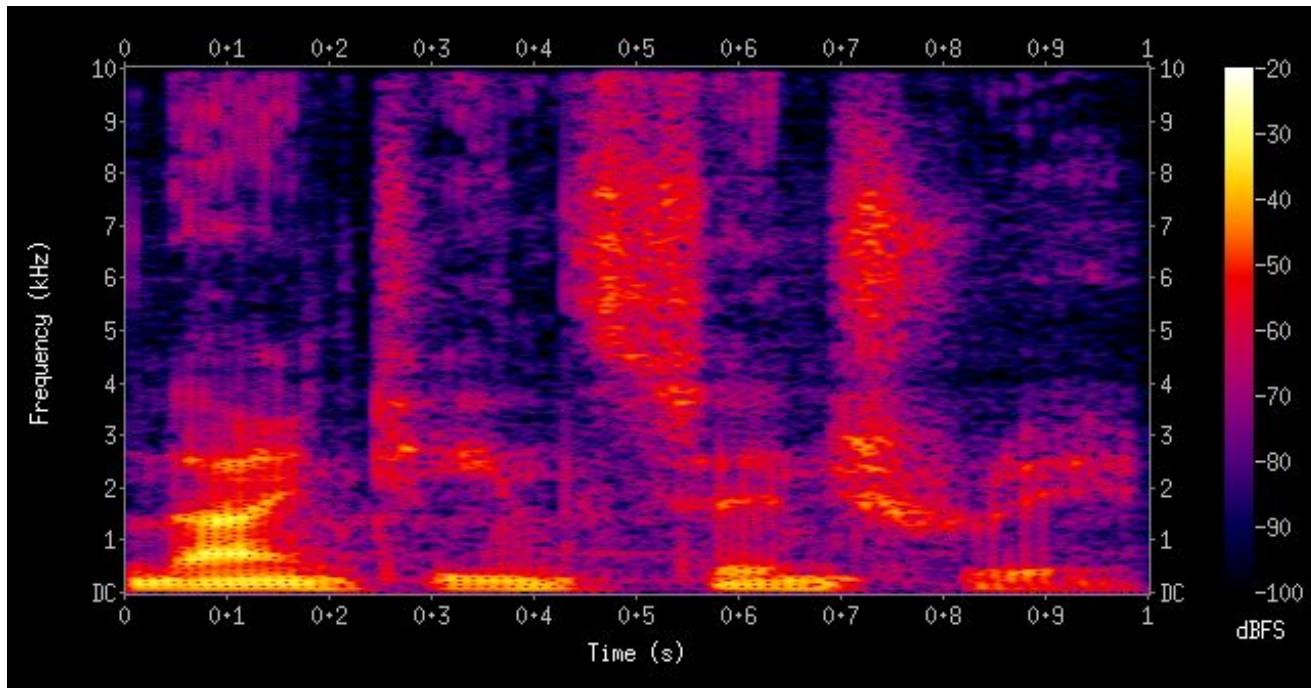


This technique is called the Short-Time Fourier Transform (STFT)!!!!!!

It solves our over-sampling problem w/out losing important information <3

We end up with a **spectrogram** that looks like this!

Short-Time Fourier Transform (STFT) Result



In our example, this is the equivalent of overlaying a couple seconds of the FTs – so we're seeing how the *frequency components change over time*

Comparative Table of Function Approximation Methods

Aspect	Fourier Transform	Taylor Series	Neural Networks
Core Idea	Represent $f(x)$ as a sum of sinusoids: $f(x) = \sum_{k=-\infty}^{\infty} c_k e^{ikx}, \quad c_k = \frac{1}{2\pi} \int f(x) e^{-ikx} dx$	Approximate $f(x)$ near a point a : $f(x) \approx \sum_{n=0}^N \frac{f^{(n)}(a)}{n!} (x - a)^n$	Learn approximation from data using nonlinear units: $f(x) \approx \sum_{i=1}^M w_i \sigma(\langle v_i, x \rangle + b_i)$ (1 hidden-layer NN)
Basis Functions	Fixed global sines/cosines.	Fixed monomials $(x - a)^n$.	Adaptive nonlinear features via activation functions.
Learning vs. Direct Computation	Coefficients c_k computed directly from integral formulas (no learning).	Coefficients $\frac{f^{(n)}(a)}{n!}$ computed directly from derivatives (no learning).	Coefficients w_i, v_i, b_i learned from data via optimization (gradient descent).
Domain of Usefulness	Oscillatory/periodic signals, frequency analysis.	Smooth analytic functions, valid locally around expansion point.	Arbitrary nonlinear, high-dimensional, non-analytic functions.
Local vs. Global	Global — each term affects all x .	Local — accurate only near a .	Both — architecture-dependent (RBFs local, deep nets capture global).
Interpretability	Coefficients \leftrightarrow frequency content.	Coefficients \leftrightarrow derivatives at expansion point.	Parameters (weights) not directly interpretable.
Convergence	Converges in L^2 for square-integrable f ; Gibbs phenomenon at discontinuities.	Converges within radius of convergence if f is analytic.	Universal Approximation Theorem: can approximate any continuous f on compact sets.
Computation	Fast Fourier Transform (FFT): $O(N \log N)$.	Easy if derivatives known; costly at high order.	Training costly (gradient descent), inference efficient.
Noise Sensitivity	Good for filtering (frequency cutoff).	Highly sensitive (derivatives amplify noise).	Can overfit to noise unless regularized.
Applications	Signal/image processing, PDEs, spectral analysis.	Physics models, perturbation methods, local expansions.	Pattern recognition, regression/classification, generative modeling.

Comparative Table of Function Approximation Methods

Aspect	Fourier Transform	Taylor Series	Neural Networks
Core Idea	Represent $f(x)$ as a sum of sinusoids: $f(x) = \sum_{k=-\infty}^{\infty} c_k e^{ikx}$, $c_k = \frac{1}{2\pi} \int f(x) e^{-ikx} dx$	Approximate $f(x)$ near a point a : $f(x) \approx \sum_{n=0}^N \frac{f^{(n)}(a)}{n!} (x-a)^n$	Learn approximation from data using nonlinear units: $f(\cdot) \approx \sum_{i=1}^M w_i \sigma(v_i \cdot x) + b$
Basis Functions	Fixed global sines/cosines.	Fixed monomials $(x-a)^n$.	Adaptive basis functions (signals).
Learning vs. Direct Computation	Coefficients c_k computed directly from integral formulas (no learning).	Coefficients $\frac{f^{(n)}(a)}{n!}$ computed directly from derivatives (no learning).	There are lots of signals where we don't really KNOW how they work... so this "math tool" is the best tool we can use
Domain of Usefulness	Oscillatory/periodic signals, frequency analysis.	Smooth analytic functions, valid locally around expansion point.	Both - able to capture global patterns.
Local vs. Global Interpretability	Speech moves through the air in waves of air pressure... so this "math tool" pretty realistically models how speech works & is cheap to use & solves our downsampling problem	accurate only near a . Coefficients \leftrightarrow derivatives at point. within radius of convergence if f is analytic.	Parameters (weights) not directly interpretable.
Convergence			Universal Approximation Theorem: can approximate any continuous f on compact sets.
Computational Cost		Easy if derivatives known; costly at high order.	Training costly (gradient descent), inference efficient.
Noise Sensitivity	Good	Highly sensitive (derivatives amplify noise).	Can overfit to noise unless regularized.
Applications	Signal/image processing, PDEs, spectral analysis.	Physics models, perturbation methods, local expansions.	Pattern recognition, regression/classification, generative modeling.

Question w/ your neighbors:

Quickly recap with your neighbor what the fourier transform does for us!

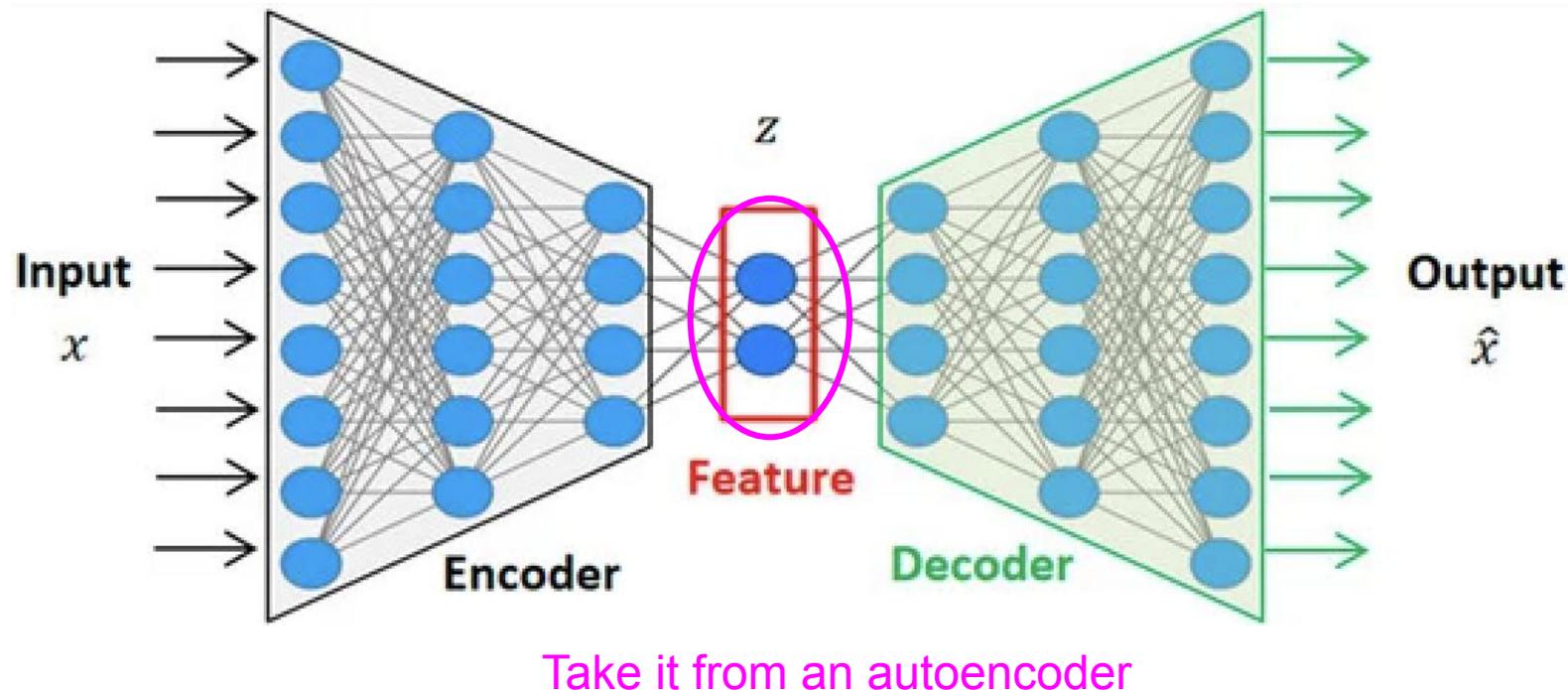
There are actually lots of ways to encode speech features –

3.1.1 Continuous Features

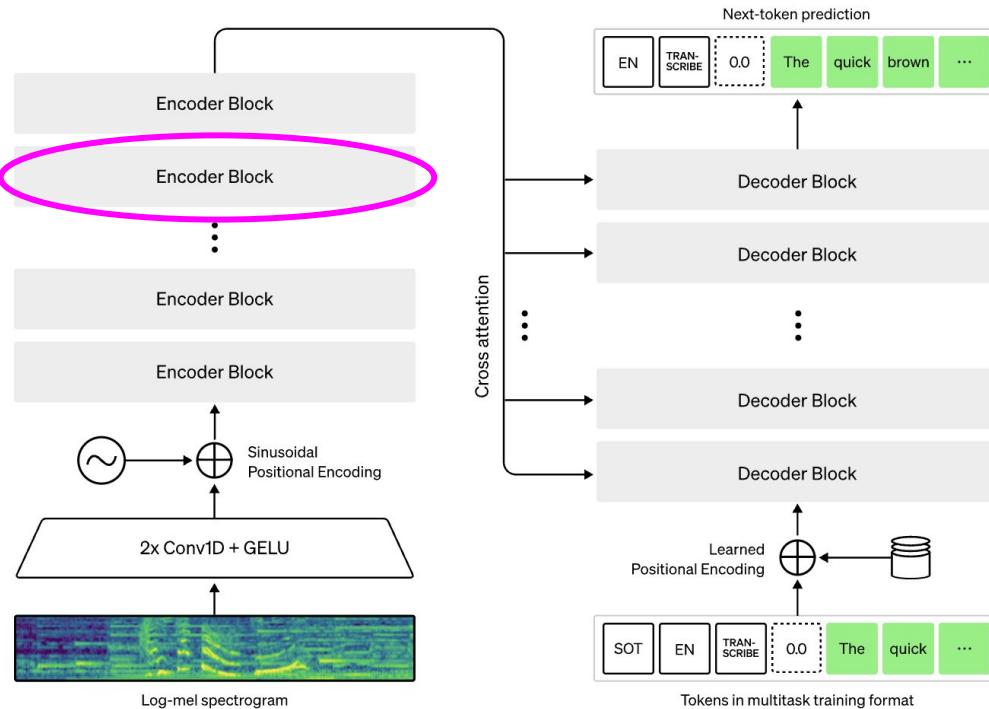
To extract informative representations from raw waveforms, a speech representation model—either a learned encoder or a digital signal processing (DSP) feature extractor—converts speech into continuous features. These continuous features may include:

-  1. Traditional spectrogram features, such as mel filter bank features (Huang et al., 2001).
- 2. Hidden representations from self-supervised learning-based (SSL) speech encoders, such as wav2vec 2.0 (Baevski et al., 2020), HuBERT (Hsu et al., 2021), or WavLM (Chen et al., 2022).
- 3. Hidden representations from supervised pre-trained models, such as Whisper (Radford et al., 2023) or USM (Zhang et al., 2023b).
- 4. Hidden representations from neural audio codec models, such as SoundStream (Zeghidour et al., 2022) or EnCodec (Défossez et al., 2023).

2. Hidden representations from self-supervised learning-based (SSL) speech encoders, such as wav2vec 2.0 (Baevski et al., 2020), HuBERT (Hsu et al., 2021), or WavLM (Chen et al., 2022).

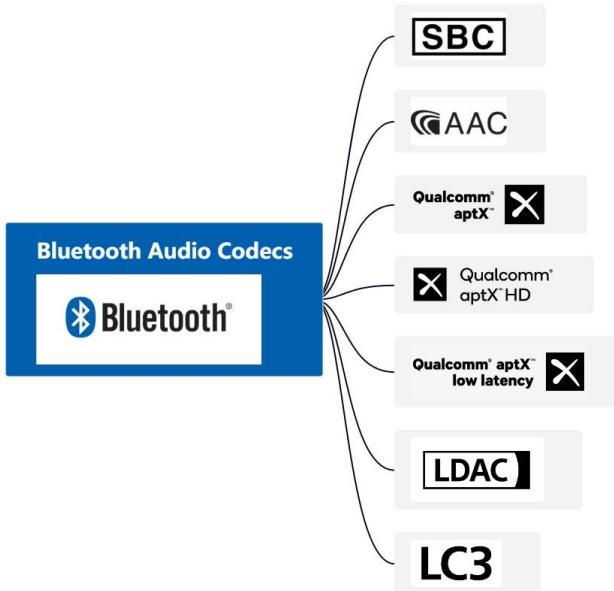


3. Hidden representations from supervised pre-trained models, such as Whisper (Radford et al., 2023) or USM (Zhang et al., 2023b).



Take it from a Whisper layer

4. Hidden representations from neural audio codec models, such as SoundStream (Zeghidour et al., 2022) or EnCodec (Défossez et al., 2023).



Codec = compressed,
lossLESS audio
representation
(big topic rn)

Question w/ your neighbors:

What are the 4 different ways
we can do a speech encoder?

Encoder

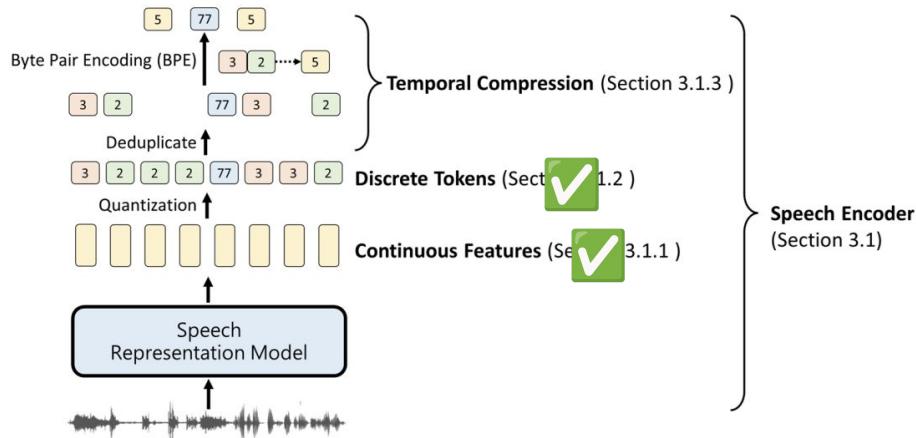


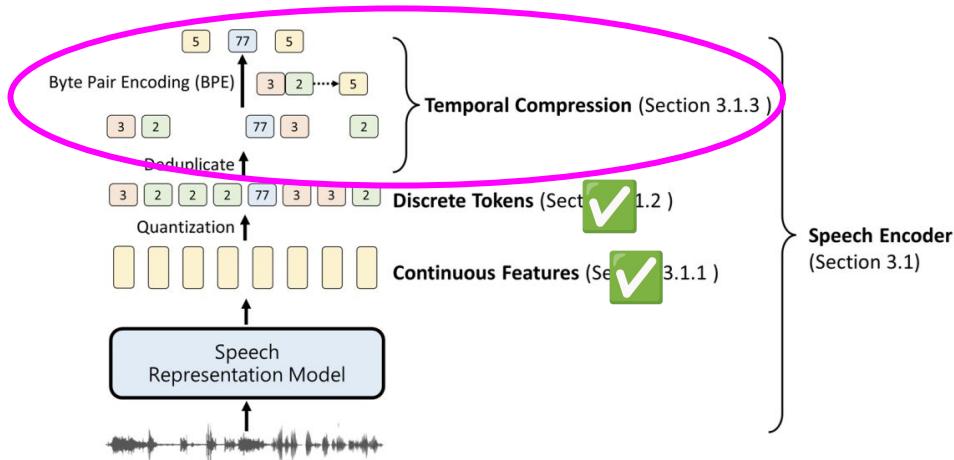
Figure 2: A general pipeline for speech encoders. Note that different encoders use different components of the pipeline. See Section 3.1 for more details.

Now we know how to go from speech signal -> tokens

(more deets in the paper)

But we're not done – there's
ANOTHER
over-sampling/compression
issue...

Encoder



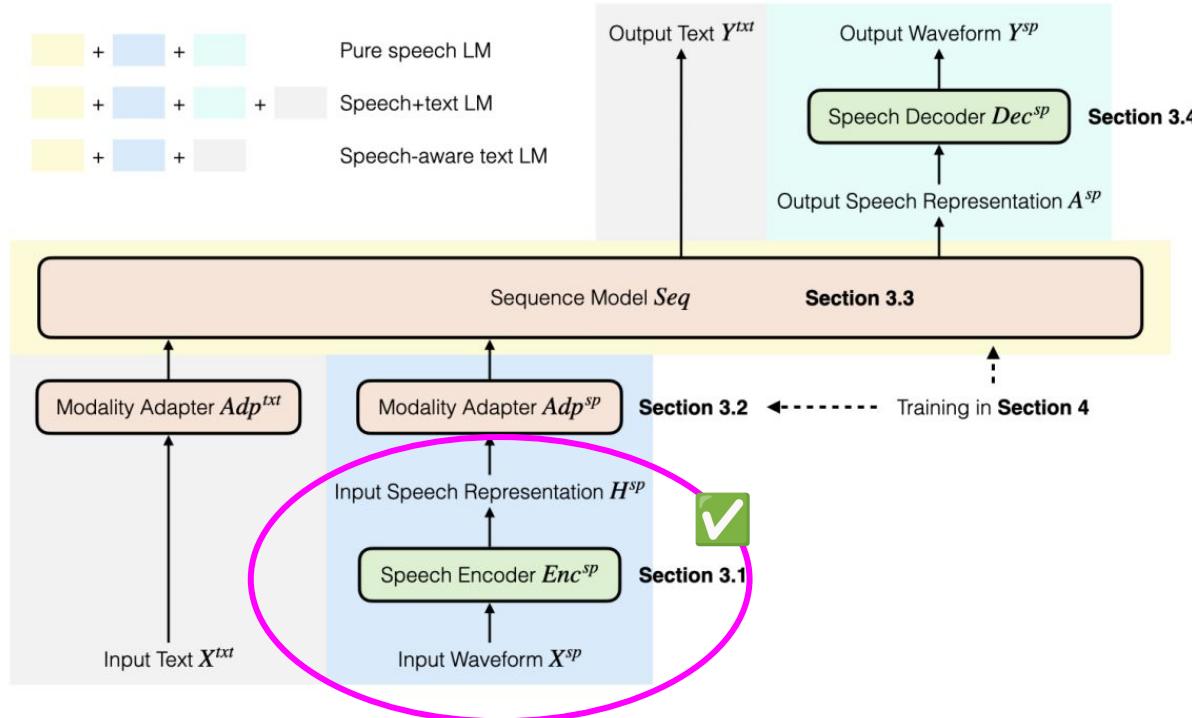
So there are lots of tools to apply here to downsample AGAIN:

- BPE
- Multi-stream capture
- Duration prediction

Figure 2: A general pipeline for speech encoders. Note that different encoders use different components of the pipeline. See Section 3.1 for more details.

Big theme: downsampling is a big deal in audio/speech world!!!!!!

Let's talk SLM architecture



Ok yay so now we know about the speech encoder

Figure 1: Overview of SLM architecture. See Sections 3 and 4 for more detailed descriptions of the components and training methods, respectively.

Let's talk SLM architecture

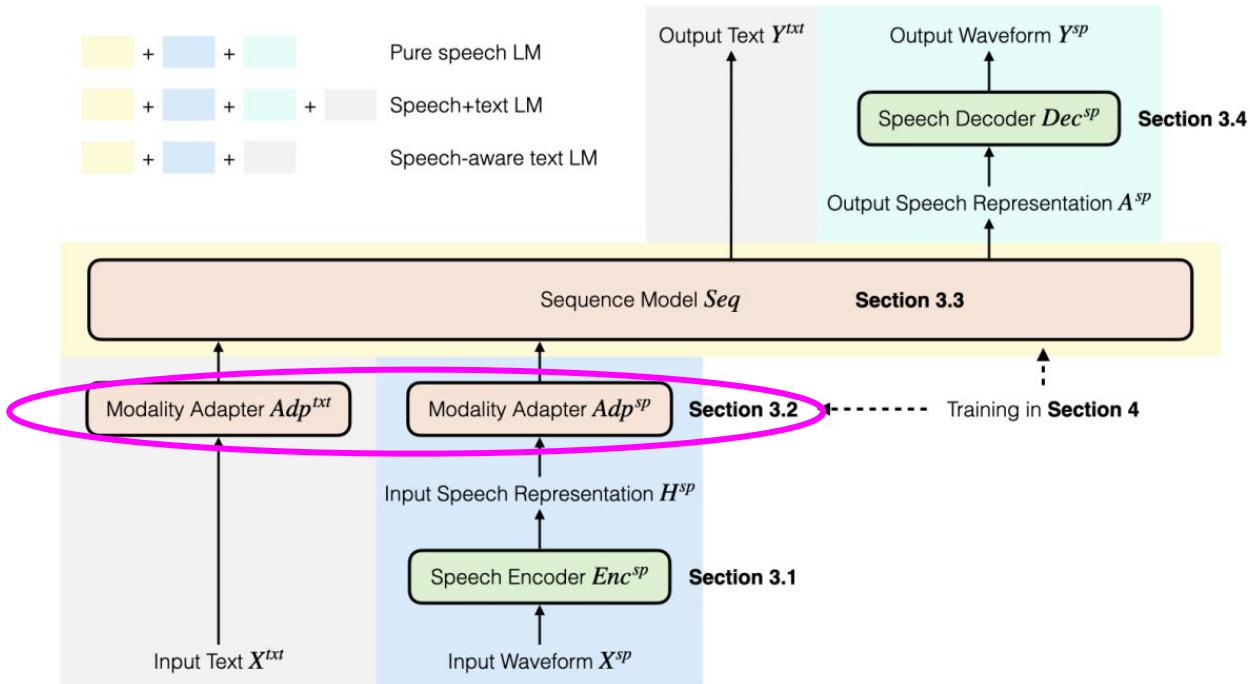


Figure 1: Overview of SLM architecture. See Sections 3 and 4 for more detailed descriptions of the components and training methods, respectively.

Let's talk about
the modality
adapters

But to really
talk about
modality
adapters we
have to talk
about **SLM**
training first!

3 Main Approaches to TRAIN SLMs

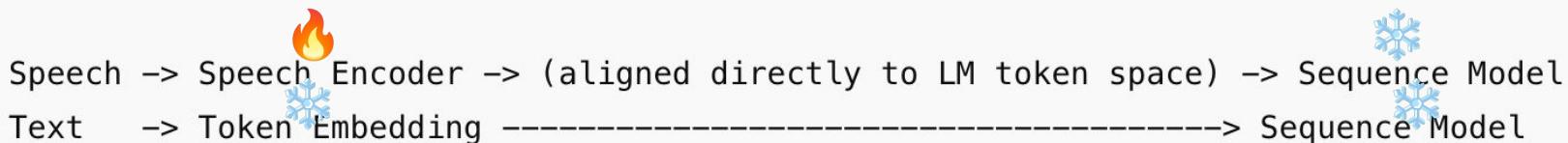
Train the (big) sequence model (aka the LLM) [most expensive \$\$\$]



Train the modality adapters [least expensive \$]



Train the speech encoder [mid-expensive \$\$] [PS this is kinda like just training 1 modality adapter]



Question w/ your neighbors:

Situation: I don't have a lot of money.

Question: Which SLM training setup should I use?

Question w/ your neighbors:

Situation: I have so much \$\$\$.

Question: Which SLM training setup should I use?

Modality Adapters

Train the modality adapters [least expensive \$]

Speech → Speech Encoder → Modality Adapter → Sequence Model
Text → Token Embedding → Modality Adapter → Sequence Model



In many SLMs (especially speech-aware text LMs), the speech encoder (Section 3.1) and the sequence model (Section 3.3) are initially developed separately and then combined. It is therefore necessary to somehow align the output of the speech encoder with the expectations of the sequence model, and this is the role of the modality adapter. The modality adapter is typically trained on downstream tasks or, in the case of speech+text LMs, as part of pre-training (see Section 4 for more details on training).

The whole point is to make the **SPEECH TOKENS and TEXT TOKENS match up correctly!!!!**

Lots of ways to do a Modality Adapter... Broad Overview

Comparative Table of Adapter Architectures in Multimodal Models

Adapter Architecture	Speech–Language	Vision–Language	Recommender–Language
Linear / MLP Projection	Project continuous speech embeddings (e.g., wav2vec2, HuBERT) into LLM token space.	Project vision encoder outputs (ViT/CNN/CLIP) into the same dimensionality as LLM embeddings.	Map item embeddings or categorical features into LLM-compatible vectors.
Quantization (Discrete Tokens)	VQ-VAE, k-means clustering to produce “pseudo-text” tokens from speech, directly fed into LLM.	Less common (though discrete VQ-image tokens exist, e.g., VQGAN), but most models use continuous projection instead.	Rare — item IDs can act as discrete tokens, sometimes directly embedded as vocab.
Cross-Attention Adapters	LLM attends to speech embeddings via cross-attention blocks (e.g., spoken QA grounding).	Popular: Flamingo inserts cross-attention layers where LLM queries image embeddings.	LLM attends over user history/item embeddings via cross-attention, aligning behavioral context with language tokens.
Prefix / Prompt Tuning	Encode speech embeddings into a sequence of pseudo-tokens prepended as “soft prompts.”	Encode image embeddings as prefix tokens before text input (BLIP-2 Q-Former, Kosmos-1).	Encode user/item history as prompt tokens describing context (P5, TALM).
Bottleneck Adapters	Small trainable modules (e.g., LSTM, Conformer bottlenecks) compress variable-length speech into manageable embeddings.	Smaller adapter layers inserted in the vision encoder or between encoder–decoder.	Trainable bottlenecks align high-dim item features with LLM hidden states.
Fusion / Multi-Modal Attention	Speech features and LLM hidden states jointly processed via fusion transformer blocks.	Multimodal transformer fusion (e.g., CLIP-LMs, PaLM-E) combining vision & text tokens.	Fusion of user/item embeddings with natural language tokens for personalized reasoning.

The whole point is to make the [SPEECH/VISION/REC] TOKENS and LLM TOKENS match up correctly!!!!

Lots of ways to do a Modality Adapter...

Connectionist Temporal Classification (CTC)-based compression: This method compresses H^{sp} (Eq. 1) according to the posterior distribution from a CTC-based speech recognizer (Gaido et al., 2021). CTC (Graves et al., 2006), a commonly used approach for ASR, assigns each time step a probability distribution over a set of label tokens, including a blank (“none of the above”) symbol. The time steps with high non-blank probabilities indicate segments that are likely to carry important linguistic information. CTC compression aggregates the frame-level labels, specifically by merging repeated non-blank labels and removing blanks. This approach produces a compressed representation intended to retain the relevant content of the original sequence while significantly reducing its length (Wu et al., 2023b; Tsunoo et al., 2024).

How CTC collapsing works



How do you
make speech
tokens & text
tokens match up
correctly?
Compression!

Lots of ways to do a Modality Adapter...

Q-Former: The Q-Former (Li et al., 2023) is an adapter that produces a fixed-length representation by encoding a speech representation sequence of arbitrary length into M embedding vectors, where M is a hyperparameter (Lu et al., 2024b).

Let the input speech representation sequence be:

$$X = \{x_1, x_2, \dots, x_{L'}\}, \quad x_i \in \mathbb{R}^{d'}, \quad (3)$$

where L' is the sequence length and d' is the dimension of the embeddings.

To achieve a fixed-length representation, Q-Former introduces M trainable query embeddings:

$$Q = \{q_1, q_2, \dots, q_M\}, \quad q_i \in \mathbb{R}^{d'}. \quad (4)$$

These queries interact with X via a cross-attention mechanism:

$$\text{Attn}(Q, X) = \text{softmax} \left(\frac{QW_Q(XW_K)^T}{\sqrt{d'}} \right) XW_V, \quad (5)$$

where W_Q , W_K , and W_V are learnable projection matrices. The result is a sequence of M embeddings.

In some approaches, instead of directly encoding the entire utterance into M vectors, a *window-level Q-Former* is applied (Yu et al., 2024; Pan et al., 2024; Tang et al., 2024) to retain temporal information. In the window-level Q-Former, the input embedding sequence is segmented, and the Q-Former is applied to each segment.

Lu et al. (2024a) compare the Q-Former with CNN-based modality adapters in a speech-aware text LM, finding that the Q-Former produces better performance on the Dynamic-SUPERB benchmark (Huang et al., 2024) (see Section 7 for more on this and other SLM benchmarks).

How do you
make speech
tokens & text
tokens match up
correctly?
Compression!

This approach
plugs more
directly into the
LLM (in the
attention
mechanism)

Lots of ways to do a Modality Adapter...

Summarizing Speech: A Comprehensive Survey → **Another great reference!** **EMNLP '25**

Fabian Retkowski¹ Maike Züfle¹ Andreas Sudmann² Dinah Pfau³
Shinji Watanabe⁴ Jan Niehues¹ Alexander Waibel^{1,4}
¹KIT ²University Bonn ³Deutsches Museum ⁴CMU

Reference	Audio Encoder	Projector	LLM
Fathullah et al. (2024)	🔥 Conformer (Gulati et al., 2020)	🔥 Linear	✳️ LLaMA-2-7B-chat (Touvron et al., 2023)
Shang et al. (2024)	🔥 Conformer (Gulati et al., 2020)	🔥 Q-Former (Li et al., 2023)	≈ LLaMA-2-7B-chat (Touvron et al., 2023)
Microsoft et al. (2025)	🔥 Conformer (Gulati et al., 2020)	🔥 MLP	✳️ Phi-4-mini-instruct (Microsoft et al., 2025)
Kang and Roy (2024)	🔥 HuBERT-Large (Hsu et al., 2021)	🔥 Linear	✳️ MiniChat-3B (Zhang et al., 2024a)
Züfle et al. (2025)	✳️ HuBERT-Large (Hsu et al., 2021)	🔥 Q-Former (Li et al., 2023)	✳️ LLaMA3.1-8B-Instruct (Grattafiori et al., 2024)
He et al. (2025)	✳️ MERaLiON-Whisper (He et al., 2025)	🔥 MLP	≈ SEA-LION V3 (He et al., 2025)
Chu et al. (2024)	🔥 Whisper-large-v3 (Radford et al., 2023)	🔥 Linear	🔥 Qwen-7B (Bai et al., 2023)
Eom et al. (2025)	✳️ Whisper-large-v2 (Radford et al., 2023)	🔥 Q-Mamba (Eom et al., 2025)	🔥 Mamba-2.8B-Zephyr (xiuyul/mamba-2.8b-zephyr)

Table 2: Overview of Audio Encoder → Projector → LLM Architectures (🔥 trainable, ✳️ frozen, ≈ LoRA)

Question w/ your neighbors:

What are some ways to set up
a modality adapter?

What is the criteria for speech/text alignment?

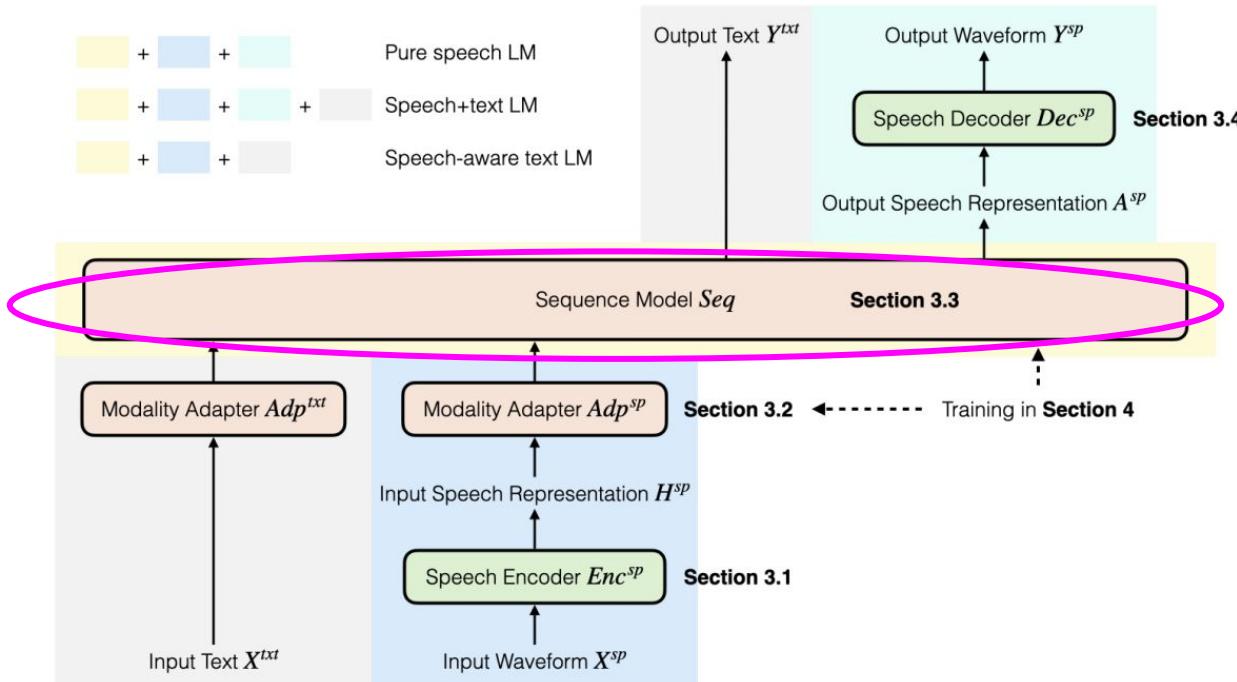
Implicit alignment Speech and text modalities can be implicitly aligned through techniques such as the “modal-invariance trick” (Fathullah et al., 2024) or behavior alignment (Wang et al., 2023a). The idea is that the model should produce identical responses regardless of the input modality, provided the input conveys the same meaning. This approach often utilizes ASR datasets. The text transcript is input to a text LLM to generate a text response, while the corresponding speech recording is input into the SLM, which is trained to generate the same text response. Another idea found to be useful for implicit alignment is training spoken LLMs for audio captioning, where a spoken LLM takes audio as input and outputs its description. It has been observed that training a spoken LLM solely through audio captioning can generalize to tasks it has never seen during training (Lu et al., 2024a;b).

Explicit alignment Speech and text modalities can also be explicitly aligned by matching speech features to corresponding text embeddings, via optimization of appropriate distance/similarity measures. For example, Wav2Prompt (Deng et al., 2024) and DiVA (Held et al., 2024) align modalities by minimizing the L_2 distance between speech features and the token embeddings of their transcripts in a text LLM while keeping the text embeddings fixed.

Question w/ your neighbors:

Why do we do speech/text alignment? What are 2 criteria?

Let's talk SLM architecture



Let's talk about
the sequence
model

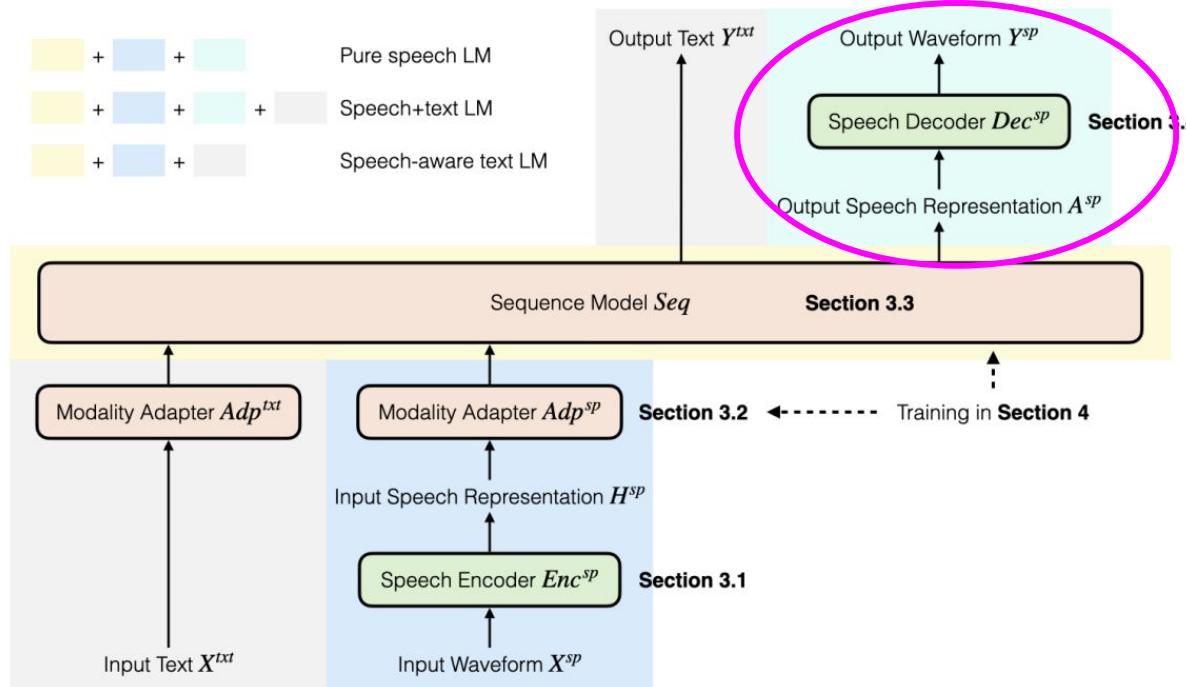
... not much to
say, this is the
core LLM

Figure 1: Overview of SLM architecture. See Sections 3 and 4 for more detailed descriptions of the components and training methods, respectively.

Question w/ your neighbors:

Tell your neighbor all the LLMs you can list off the top of your head.

Let's talk SLM architecture



Let's talk about
the **speech**
decoder

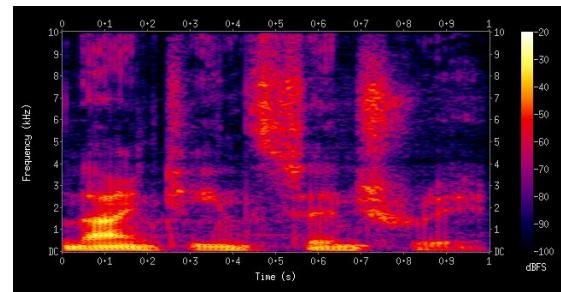
Figure 1: Overview of SLM architecture. See Sections 3 and 4 for more detailed descriptions of the components and training methods, respectively.

3.4 Speech Decoder

The speech decoder converts speech representations—whether continuous features, phonetic tokens, or audio codec tokens—**back into waveforms**. The speech decoder can take various forms:

1. Vocoder (Kong et al., 2020) for continuous features, similar to those used in traditional synthesis systems. For instance, in Spectron (Nachmani et al., 2024), a generated mel spectrogram is synthesized into audio using the WavFit vocoder (Koizumi et al., 2023).

They generate these →



2. Unit-based vocoder (Polyak et al., 2021) based on HiFi-GAN (Kong et al., 2020) for phonetic tokens. These vocoders take phonetic tokens as inputs and optionally combine them with additional information to improve synthesis quality. For example, when phonetic tokens are deduplicated, a duration modeling module is often included in the vocoder (Lakhotia et al., 2021).
3. Codec decoder (Guo et al., 2025). When the SLM generates audio codec tokens, these tokens can be input directly into the corresponding pre-trained audio neural codec decoder (without additional training) to get the waveform.

They generate codecs

Let's talk SLM architecture

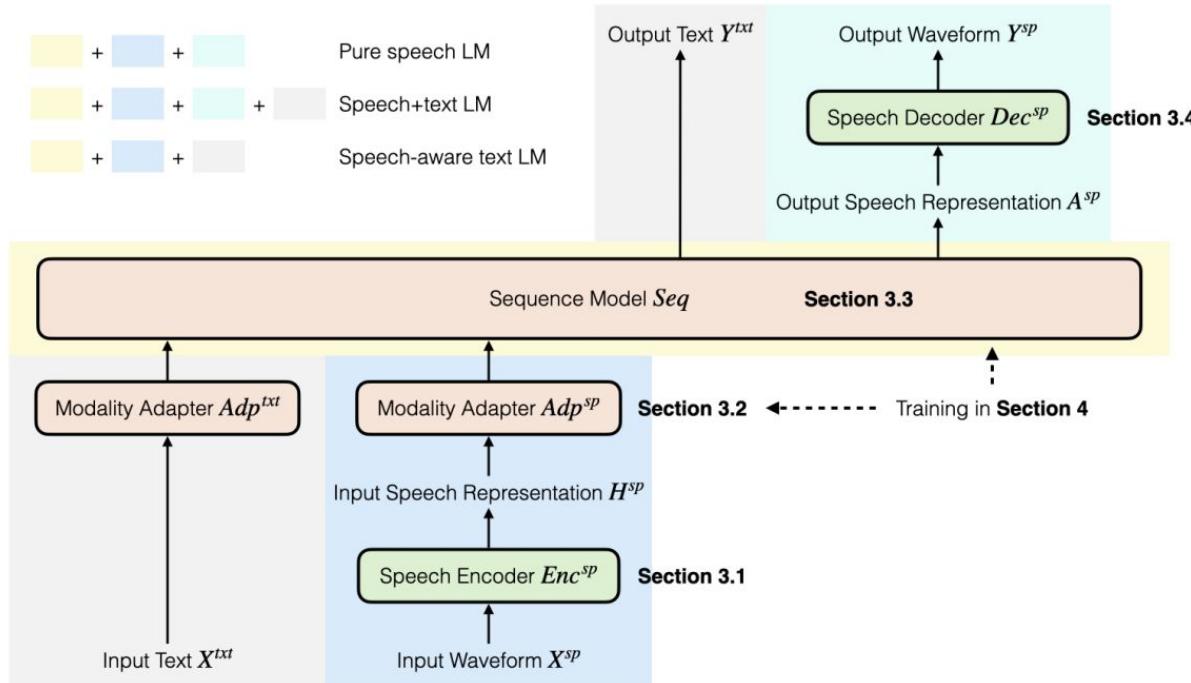


Figure 1: Overview of SLM architecture. See Sections 3 and 4 for more detailed descriptions of the components and training methods, respectively.

Ok yay!

We're done!



Question w/ your neighbors:
What does the decoder do?

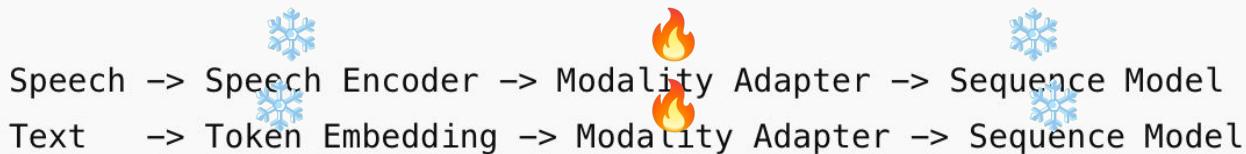
3 Main Approaches to TRAIN SLMs

Train the (big) sequence model (aka the LLM) [most expensive]

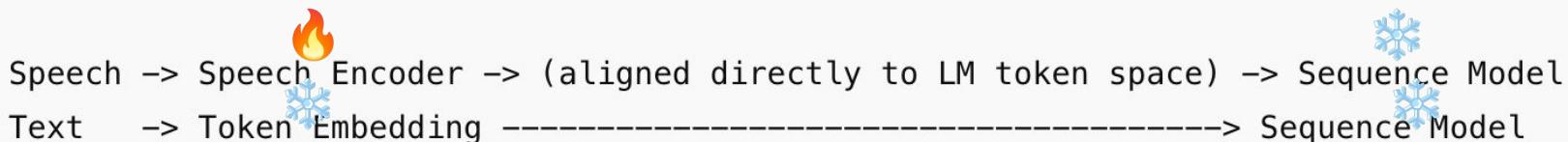


Let's revisit
how to TRAIN
SLMs!

Train the modality adapters [least expensive \$]



Train the speech encoder [mid-expensive \$\$] [PS this is kinda like just training 1 modality adapter]



They're doing the training with **special speech tokens!!!!**

$$p(\cdot | \text{speech}, \text{instruction}) \quad \left\{ \begin{array}{l} p(\cdot | \text{speech}, \text{textinstruction}), \\ p(\cdot | \text{speech}, \text{speechinstruction}) \end{array} \right. \quad \text{Either one! Just depending on setup}$$

Token / Specifier	Training Paradigm	Task	Example Input	Expected Output
<task:TRANSCRIBE>	Task-specific	ASR	[speech: "Hello, how are you today?"]	"Hello, how are you today?"
<task:SUMMARIZE>	Task-specific	Speech Summarization	[speech: lecture audio]	"The lecture explains Fourier transforms as frequency decompositions."
<task:TRANSLATE_EN>	Task-specific	Speech Translation	[speech: "Bonjour, comment ça va?"]	"Hello, how are you?"
<speech_instruction:CL ASSIFY_EMOTION>	Instruction-tuning	Emotion Detection	[speech: angry utterance]	"angry"
<speech_instruction:CO NVERSATION>	Instruction-tuning	Dialogue Response	[speech: "What time is it?"]	"It's 3:00 PM."
<speech_instruction:MU LTIMODAL_QA>	Instruction-tuning	Multimodal Fusion (Speech + Text)	[speech: "Look at this picture."] + [text: "What color is the car?"]	"The car is red."

Task-specific = "menu of fixed commands" (rigid, controlled).

Instruction tuning = "free-form instructions" (flexible, generalizable).

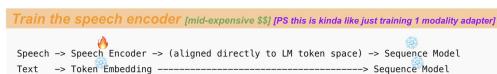
**TASK-SPECIFIC
TUNING**

**INSTRUCTION
TUNING**

Question w/ your neighbors:

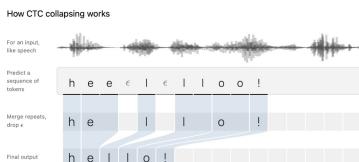
What's the point of the
special speech tokens?

So let's talk about how you **design your loss** aka objective **function for training** 🔥



TASK-SPECIFIC TUNING

INSTRUCTION TUNING



Stage	If you want to do THIS	Objective / Loss	Formula
Pretraining (AR LM)	Autoregressive next-token prediction	$\mathcal{L}_{\text{AR}} = - \sum_t \log p(y_t y_{<t}, x)$	
Adapter / Alignment	L2 alignment	$\mathcal{L}_{\text{align}} = \ f_{\text{speech}}(X^{sp}) - f_{\text{text}}(X^{txt})\ _2^2$	
	Contrastive (CLIP-style)	$\mathcal{L}_{\text{contrast}} = - \log \frac{\exp(\text{sim}(h^{sp}, h^{txt})/\tau)}{\sum_{h'} \exp(\text{sim}(h^{sp}, h')/\tau)}$	
Task-specific training	Cross-entropy w/ task token	$\mathcal{L}_{\text{task}} = - \sum_t \log p(y_t y_{<t}, X^{sp}, \langle \text{task} \rangle)$	
Instruction tuning	Cross-entropy w/ instruction	$\mathcal{L}_{\text{instr}} = - \sum_t \log p(y_t y_{<t}, X^{sp}, \text{instruction})$	
Specialized (optional)	CTC (speech alignment)	$\mathcal{L}_{\text{CTC}} = - \log p(\text{transcript} X^{sp})$	
	Reconstruction	$\mathcal{L}_{\text{recon}} = \ X^{sp} - \hat{X}^{sp}\ _2^2$	
Multi-task combo	Weighted mix	$\mathcal{L} = \sum_i \lambda_i \mathcal{L}_i$	

You'll probably end up **combo-ing** multiple of these together

Question w/ your neighbors:

Why would you combine
multiple loss functions
together?

This is cool: interruption handling w/ a duplex model

Duplex = 2 parallel streams for user and SLM, open at all times → robust to interruptions, no assumption of “turn-taking” really

Walkie-Talkie vs. Phone Call

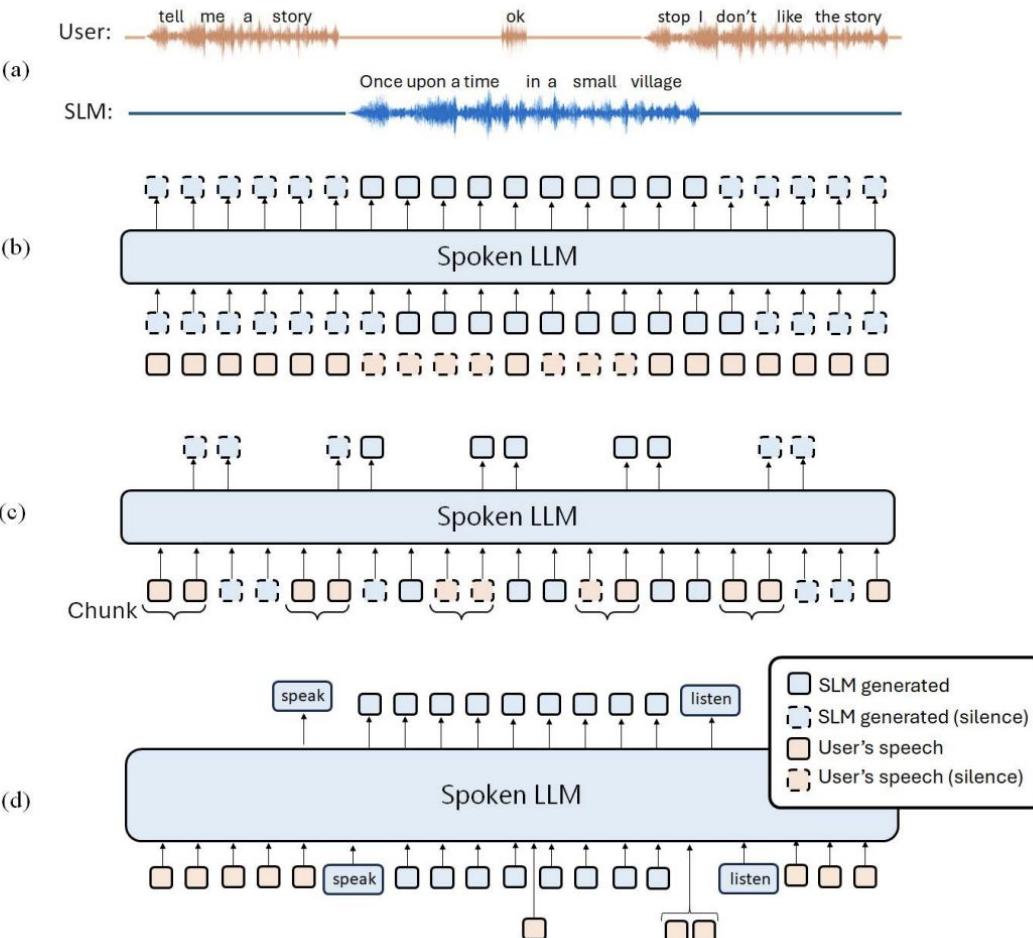


Figure 5: Full-duplex speech conversation. (a) An example of full-duplex speech conversation between a user and an SLM. (b) Dual-channel approach. (c) Time-multiplexing approach (with equal chunks). (d) Time-multiplexing approach (where the SLM controls the switching between listening and speaking modes).

8 Challenges and future work

Model architecture The optimal representation of speech within SLMs remains unclear. Speech representations in SLMs include both discrete and continuous varieties, derived from a wide range of encoders. This design choice can also influence other architectural choices in an SLM, for example depending on the information rate of the encoder and whether it encodes more phonetic or other types of information.

Another open question is determining the best method to combine speech and text, which applies to all aspects of SLM modeling and training. We have described various choices of modality adapters and approaches for interleaving speech and text. These have not been thoroughly compared, so the effect of each modeling choice is still unclear.

A final architectural challenge is that current SLMs are large and slow, making them impractical for real-time and on-device settings. To some extent this is because various compression algorithms (e.g., (Lai et al., 2021; Peng et al., 2023a; Ding et al., 2024)) and alternative architectures (e.g., Park et al. (2024)) have not been widely applied to SLMs. However, there is also an inherent efficiency challenge that arises when combining multiple pre-trained components, sometimes with different architectures and frame rates.

Finally,

**Let's talk about the
robustness of
pipeline approaches.**



WE ARE
HERE

Pipeline Approach vs. End-to-End Approach



✓ preferred by industry for *controllability* → think custom vocabulary & ease of debugging

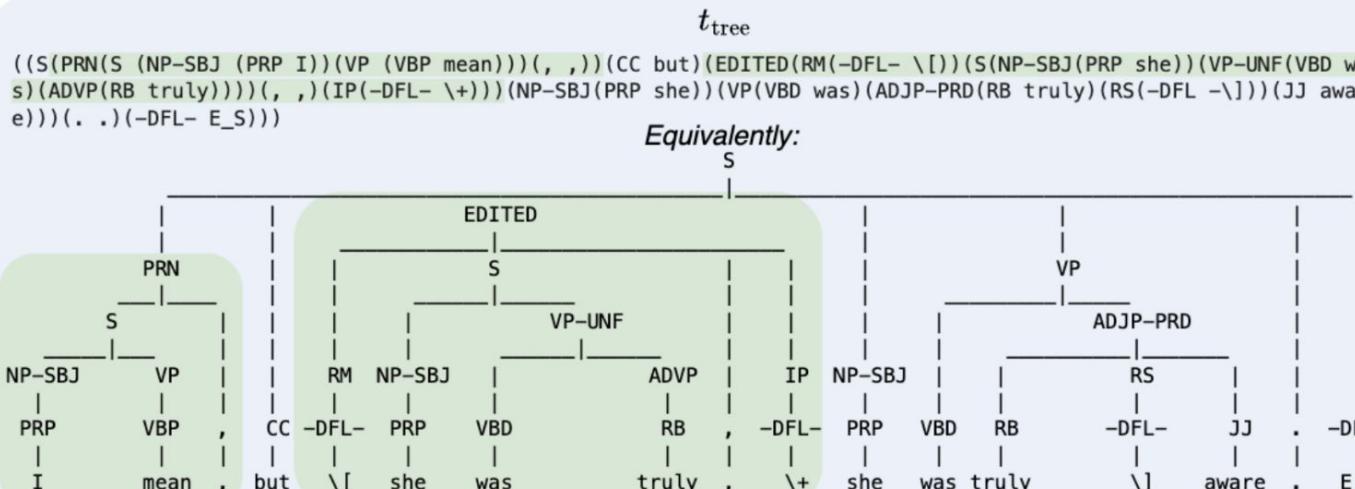
✓ currently in research stage but very promising → major barrier imo is downsampling problem (will explain)

Pipeline approaches are super susceptible to issues processing disfluencies:



Fig. 1. Removing disfluencies such as INTJ (*uh*), EDITED (*gas station* is replaced with *grocery store*), and PRN (*you know*) ensures clean text input for downstream tasks. Our

Pipeline approaches are super susceptible to issues processing disfluencies:



$t_{\text{disfluent}}$

"I mean, but she was truly, she was truly aware."

t_{tag}

[PRN, PRN, NONE, EDITED, EDITED, EDITED, NONE, NONE, NONE, NONE]

t_{fluent}

"But she was truly aware."

Quantifying the Impact of Disfluency on Spoken Content Summarization

Maria Teleki, Xiangjue Dong, James Caverlee

Texas A&M University

College Station, Texas, USA

{mariateleki, xj.dong, caverlee}@tamu.edu

LREC-COLING '24

Simple disfluencies can kill model performance.

Original

Hello and welcome to our podcast! Let's get right to it. Today we're going to be interviewing a very special guest, someone I know you guys have been excited about having on the show.

Repeats with N=3

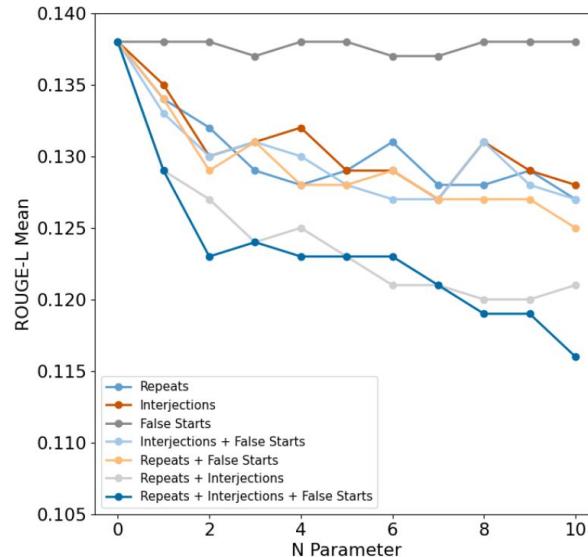
Hello and welcome to our podcast! Let's get **get get get** right to it. Today we're going to be interviewing a **a a a** very special guest, someone I know you guys have been excited about having on the show.

Interjections with N=3

Hello and welcome to our podcast! Let's get right **uh okay okay** to it. Today we're going to be interviewing a very special **um so I mean** guest, someone I know you guys have been excited about having on the show.

False Starts with N=3

Hello and welcome to our podcast! Let's get right to it. Today we're **today we're today we're today we're** going to be interviewing a very special guest, someone I know you guys have been excited about having on the show.



(a) ROUGE-L over increased N on **BART** model.

Increase N

Decrease Rouge-L

Ok but BART's old...

I agree! Here's our recent work evaluating:

- gpt-{4o,4o-mini,o4}
- llama-{1B,3B,8B,70B}
- qwen-{0.6B,1.7B,4B,8B}
- phi-4-mini
- mobileLLM{125M,350M,600M,1B}

DRES: BENCHMARKING LLMS FOR DISFLUENCY REMOVAL

Maria Teleki, Sai Janjur, Haoran Liu, Oliver Grabner, Ketan Verma, Thomas Docog, Xiangjue Dong, Lingfeng Shi, Cong Wang, Stephanie Birkelbach, Jason Kim, Yin Zhang, James Caverlee

ArXiv '25

Texas A&M University

Long story short, these models still struggle!

M	k	\mathcal{E}_P	\mathcal{E}_P	\mathcal{E}_R	Z_E	Z_I	Z_P	\parallel	\mathcal{E}_P	\mathcal{E}_P	\mathcal{E}_R	Z_E	Z_I	Z_P
gpt-4o														gpt-4o-mini
f	0	74.19 ^{+0.87}	77.91 ^{+0.63}	73.18 ^{+1.98}	83.86 ^{+0.55}	71.25 ^{+5.37}	52.52 ^{+1.91}	70.52 ^{+2.21}	75.56 ^{+0.89}	68.17 ^{+1.30}	86.13 ^{+0.39}	60.38 ^{+0.22}	55.26 ^{+0.17}	
f	1	76.12 ^{+0.63}	76.94 ^{+0.48}	78.02 ^{+0.55}	85.11 ^{+1.73}	77.37 ^{+5.64}	62.71 ^{+5.39}	68.85 ^{+0.51}	74.35 ^{+0.65}	66.02 ^{+0.44}	86.05 ^{+0.37}	57.62 ^{+0.38}	52.04 ^{+0.06}	
f	3	73.99 ^{+0.33}	81.33 ^{+0.10}	70.75 ^{+5.11}	77.83 ^{+0.30}	72.93 ^{+5.26}	45.39 ^{+5.73}	70.35 ^{+0.73}	68.70 ^{+0.62}	67.68 ^{+5.32}	66.16 ^{+0.42}	51.58 ^{+0.42}	51.58 ^{+0.42}	
f	5	73.06 ^{+0.09}	81.24 ^{+0.09}	68.75 ^{+5.35}	76.92 ^{+5.33}	71.49 ^{+4.97}	42.39 ^{+3.91}	69.03 ^{+1.17}	76.93 ^{+1.18}	65.86 ^{+4.69}	84.77 ^{+1.11}	59.46 ^{+0.51}	48.78 ^{+0.87}	
s	0	78.17 ^{+0.04}	74.48 ^{+0.83}	78.30 ^{+1.27}	82.92 ^{+0.67}	77.55 ^{+1.16}	69.80 ^{+0.56}	76.26 ^{+0.79}	75.61 ^{+0.78}	70.48 ^{+0.73}	80.20 ^{+0.23}	61.89 ^{+0.36}	65.02 ^{+0.99}	
s	1	82.38 ^{+1.11}	83.61 ^{+0.53}	81.53 ^{+1.71}	83.77 ^{+0.61}	79.74 ^{+9.44}	79.65 ^{+0.49}	77.69 ^{+0.61}	81.84 ^{+0.40}	74.34 ^{+5.81}	84.56 ^{+4.47}	65.46 ^{+0.67}	77.59 ^{+0.30}	
s	3	83.72 ^{+0.85}	84.64 ^{+0.28}	82.22 ^{+1.41}	82.42 ^{+0.90}	81.94 ^{+0.66}	78.26 ^{+2.29}	77.01 ^{+0.82}	82.10 ^{+0.53}	78.28 ^{+0.56}	64.80 ^{+0.94}	73.85 ^{+0.24}		
s	5	84.52 ^{+0.14}	88.09 ^{+0.94}	81.56 ^{+1.11}	81.15 ^{+0.70}	82.97 ^{+9.39}	77.14 ^{+0.51}	77.76 ^{+0.61}	83.31 ^{+0.42}	73.29 ^{+0.61}	83.84 ^{+0.17}	65.12 ^{+0.14}	49.69 ^{+0.04}	
medium-4o-mini														high-4o-mini
f	0	48.18 ^{+0.41}	33.04 ^{+0.81}	49.92 ^{+0.92}	84.19 ^{+0.46}	63.66 ^{+0.38}	96.51 ^{+1.17}	55.81 ^{+1.13}	41.40 ^{+0.38}	99.47 ^{+0.72}	90.49 ^{+0.32}	93.48 ^{+1.34}		
f	1	40.71 ^{+0.77}	26.19 ^{+0.57}	49.44 ^{+0.44}	84.19 ^{+0.46}	96.40 ^{+0.27}	88.64 ^{+0.44}	40.03 ^{+0.16}	31.17 ^{+0.83}	96.57 ^{+0.69}	97.42 ^{+0.44}	95.38 ^{+0.39}	97.77 ^{+0.13}	
f	3	39.65 ^{+0.77}	25.20 ^{+0.42}	49.72 ^{+0.41}	88.67 ^{+0.94}	86.84 ^{+0.51}	88.40 ^{+0.24}	42.91 ^{+0.69}	28.11 ^{+0.49}	87.27 ^{+0.69}	98.02 ^{+0.33}	96.35 ^{+0.37}	91.64 ^{+0.85}	
f	5	38.68 ^{+0.51}	24.43 ^{+0.51}	49.88 ^{+0.61}	74.45 ^{+0.01}	88.04 ^{+0.26}	41.74 ^{+0.87}	45.74 ^{+0.31}	97.84 ^{+0.51}	98.07 ^{+0.24}	97.62 ^{+0.51}	97.84 ^{+0.51}	97.07 ^{+0.04}	97.92 ^{+0.04}
Llama-1B-Instruct														Llama-3B-Instruct
f	0	35.27 ^{-7.6}	54.71 ^{+2.71}	54.72 ^{+5.13}	44.18 ^{+7.76}	35.97 ^{+7.80}	23.35 ^{+7.59}	58.45 ^{+0.06}	49.96 ^{+0.08}	78.76 ^{+0.48}	84.46 ^{+0.48}	77.40 ^{+0.91}	69.06 ^{+0.35}	
f	1	33.35 ^{+1.37}	28.67 ^{+7.55}	50.93 ^{+2.09}	75.84 ^{+0.84}	66.53 ^{+5.43}	50.45 ^{+0.33}	41.22 ^{+0.72}	81.70 ^{+0.83}	86.01 ^{+0.76}	82.32 ^{+0.14}	71.11 ^{+0.67}		
f	3	32.30 ^{+1.21}	26.65 ^{+7.91}	77.77 ^{+0.72}	70.93 ^{+0.69}	77.41 ^{+0.59}	26.44 ^{+0.20}	69.45 ^{+0.55}	76.19 ^{+0.61}	60.98 ^{+0.60}	67.52 ^{+0.61}	62.42 ^{+0.51}	45.47 ^{+0.88}	
f	5	35.60 ^{+0.37}	34.99 ^{+0.77}	68.37 ^{+0.82}	73.13 ^{+0.96}	68.51 ^{+0.72}	59.15 ^{+0.18}	48.74 ^{+0.61}	50.54 ^{+0.78}	70.08 ^{+0.51}	70.20 ^{+0.30}	77.96 ^{+0.24}	57.13 ^{+0.67}	
Llama-8B-Instruct														Llama-10B-Instruct
f	0	45.48 ^{+0.37}	31.93 ^{+0.82}	87.43 ^{+1.08}	88.57 ^{+0.81}	88.83 ^{+1.67}	61.22 ^{+2.77}	67.83 ^{+0.90}	63.90 ^{+0.75}	78.48 ^{+1.39}	81.14 ^{+0.49}	79.48 ^{+0.76}	69.48 ^{+0.23}	
f	1	37.46 ^{+1.94}	27.09 ^{+0.82}	80.38 ^{+1.29}	81.58 ^{+2.97}	75.18 ^{+0.70}	62.85 ^{+2.62}	58.99 ^{+0.77}	74.55 ^{+0.56}	74.84 ^{+0.50}	74.84 ^{+0.44}	76.44 ^{+0.50}	66.90 ^{+0.51}	
f	3	30.32 ^{+0.81}	18.02 ^{+1.76}	99.94 ^{+0.60}	100.00 ^{+0.00}	99.94 ^{+0.39}	99.85 ^{+0.01}	58.37 ^{+0.72}	82.89 ^{+0.47}	85.32 ^{+0.88}	84.17 ^{+0.30}	75.74 ^{+0.87}		
f	5	30.33 ^{+0.95}	18.02 ^{+1.76}	100.00 ^{+0.00}	100.00 ^{+0.00}	100.00 ^{+0.00}	100.00 ^{+0.00}	53.37 ^{+0.72}	44.39 ^{+0.78}	82.87 ^{+0.47}	83.37 ^{+0.82}	85.37 ^{+0.73}	53.63 ^{+0.39}	34.55 ^{+0.02}
Qwen-3B-struct														Qwen-3B-struct
f	0	18.04 ^{+0.01}	43.29 ^{+0.58}	16.04 ^{+0.61}	22.51 ^{+2.79}	14.26 ^{+4.23}	11.56 ^{+0.36}	10.25 ^{+0.49}	86.02 ^{+0.04}	5.91 ^{+0.01}	10.36 ^{+0.09}	4.10 ^{+0.26}	2.75 ^{+0.04}	
f	1	22.36 ^{+0.02}	29.02 ^{+0.34}	35.09 ^{+0.65}	28.30 ^{+3.63}	27.18 ^{+5.02}	10.07 ^{+0.21}	79.60 ^{+0.52}	13.39 ^{+0.55}	16.82 ^{+0.79}	11.84 ^{+0.26}	10.70 ^{+0.97}		
f	3	21.17 ^{+0.28}	38.69 ^{+1.38}	20.34 ^{+2.44}	24.92 ^{+7.04}	19.29 ^{+3.20}	14.89 ^{+0.91}	7.78 ^{+0.20}	88.20 ^{+0.71}	5.09 ^{+0.14}	8.82 ^{+0.30}	3.24 ^{+0.28}	2.98 ^{+0.14}	
f	5	19.87 ^{+0.52}	40.95 ^{+0.37}	19.67 ^{+0.41}	24.09 ^{+1.81}	17.87 ^{+5.44}	16.06 ^{+0.62}	8.64 ^{+0.42}	84.52 ^{+0.85}	6.66 ^{+0.22}	10.71 ^{+0.12}	4.88 ^{+0.15}	4.17 ^{+0.08}	
Qwen-3B														Qwen-3B
s	0	43.74 ^{+1.31}	44.90 ^{+0.86}	44.05 ^{+0.56}	80.78 ^{+1.71}	33.94 ^{+0.34}	11.94 ^{+0.35}	39.00 ^{+0.24}	74.11 ^{+0.04}	24.79 ^{+0.04}	39.19 ^{+1.40}	13.99 ^{+0.61}	14.95 ^{+0.02}	
s	1	51.97 ^{+0.69}	47.23 ^{+0.23}	59.72 ^{+0.79}	70.39 ^{+0.95}	55.32 ^{+0.35}	53.93 ^{+0.38}	35.78 ^{+0.25}	65.77 ^{+0.19}	25.10 ^{+0.04}	33.94 ^{+0.64}	6.18 ^{+0.45}		
s	3	48.90 ^{+0.67}	54.09 ^{+0.84}	50.04 ^{+1.14}	62.76 ^{+0.75}	52.82 ^{+2.4}	42.36 ^{+0.46}	31.20 ^{+0.07}	67.20 ^{+0.04}	18.55 ^{+0.01}	24.86 ^{+0.37}	5.09 ^{+0.31}		
s	5	48.38 ^{+0.33}	48.74 ^{+0.73}	50.29 ^{+0.45}	58.66 ^{+0.18}	50.05 ^{+0.36}	36.34 ^{+0.73}	34.46 ^{+0.45}	81.04 ^{+0.01}	22.31 ^{+0.01}	23.00 ^{+0.04}	29.00 ^{+0.07}	5.08 ^{+0.01}	
Qwen-3B-struct														Phi-4-mini-struct
f	0	66.39 ^{+0.56}	75.58 ^{+6.33}	64.55 ^{+0.98}	62.47 ^{+4.32}	70.30 ^{+2.66}	49.66 ^{+0.16}	71.04 ^{+0.77}	69.94 ^{+12.12}	76.17 ^{+5.49}	75.24 ^{+8.17}	79.32 ^{+0.48}	67.37 ^{+0.01}	
f	1	51.89 ^{+0.69}	59.63 ^{+0.49}	54.22 ^{+0.68}	55.44 ^{+3.33}	59.91 ^{+4.54}	63.76 ^{+4.43}	68.86 ^{+3.06}	54.54 ^{+5.36}	69.69 ^{+7.63}	67.32 ^{+6.07}	54.80 ^{+8.83}		
f	3	62.29 ^{+0.67}	80.58 ^{+4.28}	57.13 ^{+0.73}	54.85 ^{+5.03}	65.16 ^{+0.98}	39.32 ^{+0.83}	71.77 ^{+0.73}	76.82 ^{+0.73}	70.78 ^{+0.65}	68.89 ^{+3.34}	77.84 ^{+0.83}	54.49 ^{+0.83}	
f	5	68.45 ^{+0.55}	67.96 ^{+9.44}	65.05 ^{+2.29}	70.68 ^{+0.41}	62.02 ^{+0.69}	71.46 ^{+0.71}	78.32 ^{+0.63}	66.17 ^{+0.77}	77.70 ^{+0.93}	69.37 ^{+2.39}	62.99 ^{+2.24}	59.76 ^{+0.96}	
Phi-4-m mini-struct														Phi-4-m mini-struct
s	0	33.23 ^{+0.01}	20.58 ^{+0.77}	91.42 ^{+1.51}	91.95 ^{+0.72}	83.76 ^{+2.10}	34.25 ^{+0.21}	21.85 ^{+0.15}	82.74 ^{+0.67}	82.74 ^{+0.67}	80.86 ^{+0.86}	84.47 ^{+0.18}	85.25 ^{+0.16}	
s	1	31.90 ^{+0.16}	21.13 ^{+0.87}	89.38 ^{+1.02}	92.02 ^{+0.24}	90.79 ^{+0.83}	35.16 ^{+0.24}	22.20 ^{+0.03}	80.86 ^{+0.69}	84.17 ^{+0.69}	86.45 ^{+0.39}	86.45 ^{+0.39}		
s	3	34.40 ^{+0.16}	22.30 ^{+0.32}	80.11 ^{+0.84}	84.82 ^{+0.12}	70.45 ^{+0.57}	36.35 ^{+0.33}	28.03 ^{+0.09}	54.38 ^{+0.23}	60.63 ^{+0.27}	54.72 ^{+0.83}	54.72 ^{+0.83}	54.07 ^{+0.09}	
s	5	34.99 ^{+0.17}	22.64 ^{+1.11}	80.40 ^{+0.84}	83.64 ^{+0.29}	84.12 ^{+0.32}	66.45 ^{+0.32}	37.84 ^{+0.21}	64.55 ^{+0.22}	67.10 ^{+0.22}	73.07 ^{+0.24}	75.58 ^{+0.31}	62.29 ^{+0.12}	
Phi-4-mLLM-125M														Phi-4-mLLM-350M
s	0	33.23 ^{+0.01}	20.58 ^{+0.77}	91.42 ^{+1.51}	91.95 ^{+0.72}	83.76 ^{+2.10}	34.25 ^{+0.21}	21.85 ^{+0.15}	82.74 ^{+0.67}	82.74 ^{+0.67}	80.86 ^{+0.86}	84.47 ^{+0.18}	85.25 ^{+0.16}	
s	1	36.13 ^{+0.16}	27.10 ^{+0.88}	56.91 ^{+0.41}	61.45 ^{+2.22}	55.20 ^{+0.36}	34.25 ^{+0.21}	22.20 ^{+0.03}	82.11 ^{+0.69}	63.68 ^{+0.27}	66.74 ^{+0.72}	63.90 ^{+0.55}	57.95 ^{+0.36}	
s	3	36.18 ^{+0.44}	29.21 ^{+0.98}	49.84 ^{+0.57}	57.54 ^{+2.21}	48.52 ^{+0.63}	40.83 ^{+0.75}	33.55 ^{+0.56}	23.94 ^{+0.71}	59.55 ^{+0.39}	61.97 ^{+0.41}	61.06 ^{+0.38}	52.16 ^{+0.97}	
s	5	36.98 ^{+0.44}	30.83 ^{+0.65}	34.74 ^{+0.99}	64.79 ^{+0.32}	69.29 ^{+0.21}	48.86 ^{+0.34}	37.15 ^{+0.36}	25.41 ^{+0.71}	72.40 ^{+0.36}	74.39 ^{+0.97}	76.00 ^{+0.41}	59.47 ^{+0.67</}	

So we make a set of concrete recommendations to help these models perform better on disfluent, spoken data!

3. RESULTS & RECOMMENDATIONS

We analyze the disfluency removal behavior of LLMs and provide recommendations (R1-R9).

Open-Source vs. Proprietary. Looking to Table 3, proprietary models (gpt-4o, gpt-4o-mini) achieve the highest scores, with margins of 10–15 points over the best open-source alternatives. We attribute this to training exposure to Whisper-transcribed speech data [24]. **(R1) Proprietary models are currently the most reliable for production systems, while open-source models require targeted augmentation with spoken data.**

Segmentation (s) vs. Full Input (f). Segmenting transcripts consistently improves both mean performance and stability, e.g., gpt-4o improves from $\mathcal{E}_F=76.13$ (f) to 82.38 (s) at $k=1$. This supports prior evidence of long-context degradation in LLMs [25, 26]. **(R2) Segmentation is an effective preprocessing step that should be applied.**

Few-Shot Sensitivity (k). Increasing k does not uniformly improve results. Small models (e.g., MobileLLM) gain slightly, but others show degradation (e.g. Llama-3B/8B/70B) when more examples are provided. **(R3) Few-shot prompting should be used with caution, as some model families misinterpret exemplars and over-edit fluent text.**

Disfluency Category Performance. Z -Scores show that EDITED nodes are handled well, but INTJ and PRN nodes are frequently missed, despite prior work suggesting these are the easiest to detect [19, 17]. **(R4) Future modeling should focus on under-served categories (INTJ, PRN) to improve robustness across all disfluency types.**

Over-Deletion Failures. Several models (e.g., Llama-8B, o4-mini) achieve near perfect recall but at the cost of very low precision, deleting fluent tokens. Segmentation often mitigates this collapse mode. **(R5) Segment-level evaluation**

helps reduce over-deletion risk.

Under-Deletion Failures. Some models (e.g., Qwen series) exhibit the opposite trend of over-deletion, achieving high precision but low recall (purple). These models preserve most fluent tokens but fail to remove many true disfluencies, especially in INTJ and PRN categories. This reflects conservative editing strategies and limited exposure to conversational disfluency distributions. **(R6) Models prone to under-deletion require additional filtering or targeted fine-tuning to ensure sufficient disfluency coverage.**

Reasoning-Oriented Models. Models tuned for reasoning (o4-mini, Phi-4) perform poorly, showing high recall but extreme over-deletion (blue). **(R7) Reasoning capability does not translate to disfluency removal; specialized evaluation remains necessary.**

Impact of Model Size. Model scaling generally improves disfluency removal, with Qwen, GPT, and Llama families showing upward trends. However, gains are nonlinear – e.g., Qwen3-1.7B underperforms both smaller and larger variants – likely due to training data or optimization differences rather than capacity limits. **(R8) Model choice should be guided by empirical benchmarks on target domains and disfluency categories rather than size alone.**

Fine-Tuning and Generalization. Looking to Tables 4 and 5, fine-tuning improves performance to near SOTA levels (e.g., gpt-4o-mini_{ft} achieves $\mathcal{E}_P=96.6$), but evaluation on GSM8K, MMLU, and CoQA shows degraded performance on unrelated tasks. **(R9) Fine-tuning is suitable for dedicated disfluency pipelines, but not for general-purpose conversational models.**

What about ASR systems? Are they also bad at processing disfluencies?

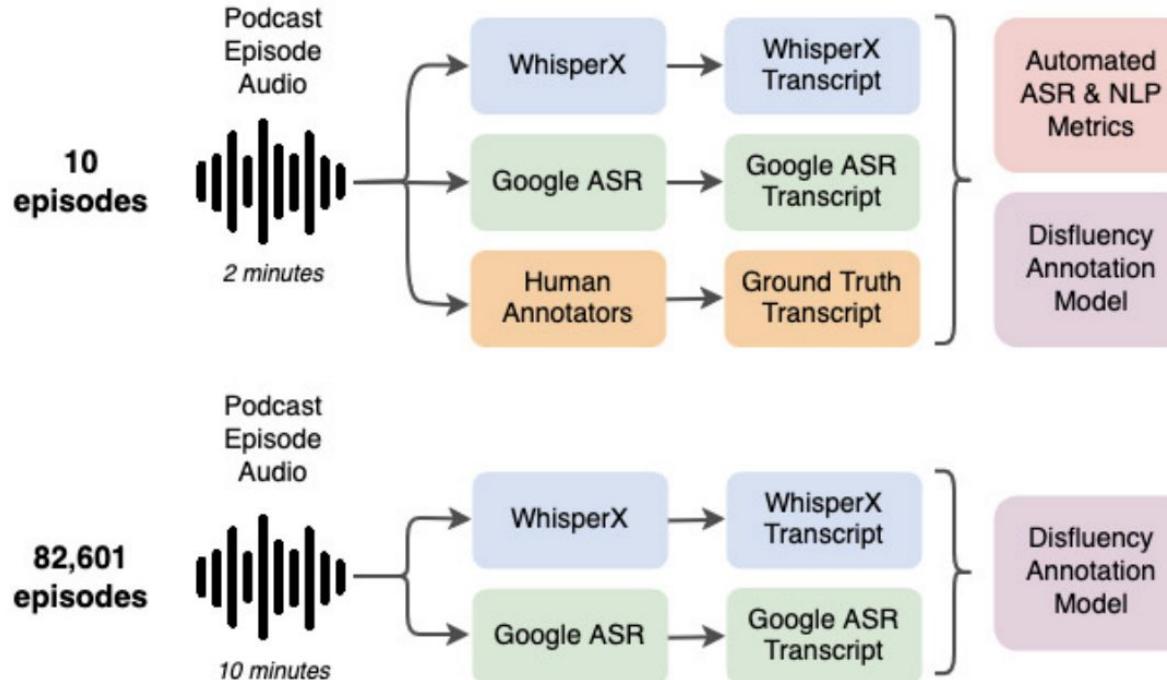
... Yes

Comparing ASR Systems in the Context of Speech Disfluencies
INTERSPEECH '24

Maria Teleki¹, Xiangjue Dong¹, Soohwan Kim¹, James Caverlee¹

¹Texas A&M University, USA

{mariateleki, xj.dong, cocomox26, caverlee}@tamu.edu



Question w/ your neighbors:

Why do you think LLMs are so
BAD at dealing with
disfluencies?

Collaborators



James Caverlee
(my advisor)



Xiangjue Dong



Haoran Liu



Sai Tejas
Janjur



Thomas
Docog



Oliver Grabner



Soohwan Kim



Stephanie
Birkelbach



Rohan
Chaudhury



Ketan Verma



Chengkai Liu



Yin Zhang

+ Jason Kim, Lingfeng Shi, Cong Wang, and shoutout to work-in-progress collaborators too :)

Conversational AI

MARIA TELEKI

Slides are
posted!



Texas A&M University
Department of Computer Science & Engineering

