



thesis defense presentation

By Amir Djanibekov

MSc Student at NLP department

INTRODUCTION

The image displays two side-by-side screenshots illustrating the search for large language models.

Hugging Face Model Hub: The left screenshot shows the Hugging Face Model Hub interface. A search bar at the top contains the query "large language model". Below it, a "Models" section lists 263 results. The first few items include "h2oai/h2o-danube2-1.8b-chat", "h2oai/h2o-danube2-1.8b-sft", "fbellame/llama2-pdf-to-quizz-13b", and "h2oai/h2o-danube-1.8b-chat". The sidebar on the left provides filtering options for tasks like "text-generation-inference" and "Inference Endpoints".

GitHub Search Results: The right screenshot shows GitHub search results for "large language models". A search bar at the top contains the same query. The results are filtered by "Code" and show 6k results. The top result is "BradyFU/Awesome-Multimodal-Large-Language-Models", which is described as "Latest Papers and Datasets on Multimodal Large Language Models, and Their Evaluation." Other results include "databricks-academy/large-language-models" and "MuhammadMoinFaisal/LargeLanguageModelsProjects".

<https://huggingface.co/models?other=large+language+model&sort=trending> <https://github.com/search?q=large+language+models&type=repositories>

INTRODUCTION

....public models that you can experiment with

....ongoing/completed works with large language models

The screenshot shows the Hugging Face website interface. At the top, there's a navigation bar with 'Hugging Face' logo, a search bar, and tabs for 'Models', 'Datasets', and 'Other'. The 'Other' tab is selected, showing a count of 263. Below this, there's a filter bar with 'large language model' selected. The main area displays a list of public models, each with a thumbnail, name, description, and statistics like 'Text Generation', 'Updated 3 days ago', and '343' likes.

<https://huggingface.co/models?other=large+language+model&sort=trending>

The screenshot shows a GitHub search results page for 'large language models'. A red arrow points from the 'Models' section of the Hugging Face screenshot to this GitHub search. The search bar at the top shows 'large language models'. The results are filtered by 'Code' and show 6k results. The first result is a repository named 'BradyFU/Awesome-Multimodal-Large-Language-Models' with a description of 'Latest Papers and Datasets on Multimodal Large Language Models, and Their Evaluation.' Other results include 'databricks-academy/large-language-models' and 'MuhammadMoinFaisal/LargeLanguageModelsProjects'.

<https://github.com/search?q=large+language+models&type=repositories>

INTRODUCTION

....public models that you can experiment with

....ongoing/completed works with large language models

The screenshot shows the Hugging Face website interface. At the top, there is a navigation bar with links for Tasks, Libraries, Datasets, Languages, Licenses, and Other (which has 1 item). Below this is a search bar labeled "Search models, datasets, users...". A red arrow points from the text "....public models that you can experiment with" to the "Models" tab, which is highlighted with a red circle and shows the number 263. The main content area displays several model cards, each with a thumbnail, name, description, and metrics like text generation and updates.

Models 263

- h2oai/h2o-danube2-1.8b-chat
- h2oai/h2o-danube2-1.8b-sft
- fbellame/llama2-pdf-to-quizz-13b
- h2oai/h2o-danube-1.8b-chat

The screenshot shows a GitHub search results page for "large language models". A red arrow points from the text "....ongoing/completed works with large language models" to the search bar, which also has "large language models" typed into it. The results show a list of repositories, with the first one being "BradyFU/Awesome-Multimodal-Large-Language-Models" and having 6k results. A red circle highlights the "6k results (212 ms)" text.

large language models

Filter by

- Code 2.6M
- Repositories 6k
- Issues 24k
- Pull requests 108k
- Discussions 1k
- Users 280
- Commits 33k
- Packages 9
- Wikis 8k
- Topics 11
- Marketplace 0

6k results (212 ms)

- BradyFU/Awesome-Multimodal-Large-Language-Models
- databricks-academy/large-language-models
- MuhammadMoinFaisal/LargeLanguageModelsProjects

<https://huggingface.co/models?other=large+language+model&sort=trending>

<https://github.com/search?q=large+language+models&type=repositories>

pre-train.... fine-tune.... instruction-tune....



MOTIVATION

The image shows two side-by-side GitHub search results pages. The left page is for 'vision large language models' and the right page is for 'speech large language model'. Both results are filtered by 'Repositories'. The 'vision' search yields 121 results, while the 'speech' search yields 51 results. Arrows point from the GitHub URLs below each search interface up to their respective search results.

vision large language models

121 results (199 ms)

PKU-YuanGroup/MoE-LLaVA
Mixture-of-Experts for Large Vision-Language Models
moe multi-modal mixture-of-experts large-vision-language-model
Python · 1.6k · Updated 23 days ago

Ucas-HaoranWei/Vary
Official code implementation of Vary: Scaling Up the Vision Vocabulary of Large Vision L
Python · 1.5k · Updated 9 days ago

PacktPublishing/Pretrain-Vision-and-Large-Language-Models-in-Python
Pretrain Vision and Large Language Models in Python, Published by Packt
Jupyter Notebook · 73 · Updated on Dec 22, 2023

Filter by

- Code
- Repositories**
- Issues
- Pull requests
- Discussions
- Users
- Commits
- Packages
- Wikis
- Topics
- Marketplace

Languages

speech large language model

51 results (270 ms)

Onutation/SpeechGPT
SpeechGPT Series: Speech Large Language Models
Python · 873 · Updated 11 days ago

yI4579/StyleTTS2
StyleTTS 2: Towards Human-Level Text-to-Speech through Style Diffusion and Adversarial Training with Large Speech Language Models
text-to-speech deep-learning pytorch tts speech-synthesis
Python · 4k · Updated 16 days ago

Onutation/USLM
Unified Speech Language Model for paper "SpeechTokenizer: Unified Speech Tokenizer for Speech Large Language Models"(ICLR 2024)
Python · 104 · Updated on Sep 14, 2023

YuanGongND/ltu
Code, Dataset, and Pretrained Models for Audio and Speech Large Language Model "Listen, Think, and Understand".
audio deep-learning speech-recognition audio-processing large-language-models
Python · 272 · Updated 23 days ago

Filter by

- Code
- Repositories**
- Issues
- Pull requests
- Discussions
- Users
- More

Languages

- Python
- Jupyter Notebook
- C++
- HTML
- Java
- JavaScript
- Kotlin
- More languages...

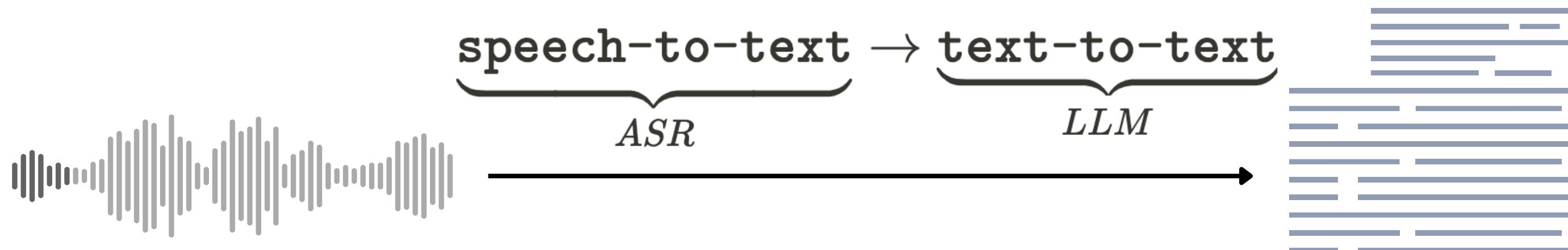
<https://github.com/search?q=vision+large+language+model&type=repositories>

<https://github.com/search?q=speech+large+language+model&type=repositories>

vision-language applications > speech-language applications

.....

MOTIVATION



Speech recognition is a sequence classification problem



h	h	h	h	h	h	h	h	h
e	e	e	e	e	e	e	e	e
l	l	l	l	l	l	l	l	l
o	o	o	o	o	o	o	o	o
€	€	€	€	€	€	€	€	€

h e € l l € l l o o
h h e l l € € l € o o
€ e € l l € € l o o

h e l l o
e l l o
h e l o

.... usually done by CTC training objective

$$\underbrace{B}_{\text{Batch size}} \times \underbrace{T}_{\text{Temporal dimension}} \times \underbrace{D}_{\text{Hidden dimension}} \longrightarrow \underbrace{B}_{\text{Batch size}} \times \underbrace{T}_{\text{Temporal dimension}}$$

.... index from vocab

OVERALL

- LLMs approximate human-like understanding
 - integrating other modalities enriches LLMs contextual information
- Vision-language modeling already reached outstanding results in image captioning
- Discriminative Speech Models are good but Deep Discriminative Speech Models might better
 - Decoding tokens to text at decoder side
- (Motivation) Text Generative Speech Models another promising direction
 - LLM integration as vision-language

LIMITATIONS

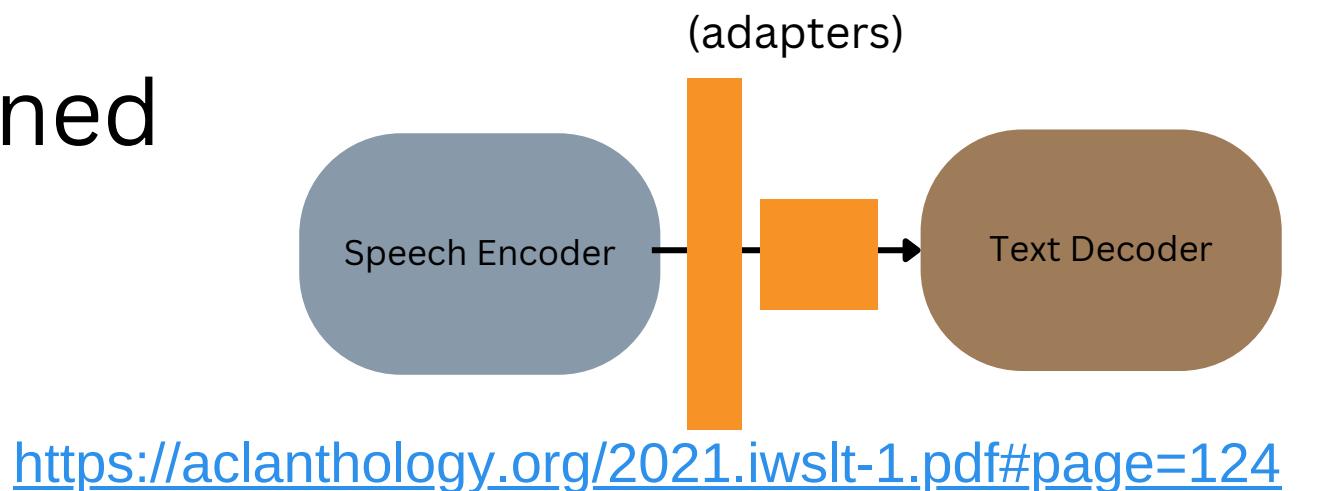
- Training LLM is compute exhaustive and expensive procedure
 - x2 pain collecting data
 - x3 pain verifying collected data
 - x4 pain training
 - x100 pain evaluating
- Speech-language modeling with LLM
 - next ‘x’ pain?

LIMITATIONS

- Training LLM is compute exhaustive and expensive procedure
 - x2 pain collecting data
 - x3 pain verifying collected data
 - x4 pain training
 - x100 pain evaluating
- Speech-language modeling with LLM
 - next ‘x’ pain?
 - we’ll try parameter efficient way by freezing LLM and speech encoder while training only modality connector/adapter

RELATED WORKS

- End-to-End Speech Translation with Pre-trained Models and Adapters: UPC at IWSLT 2021



<https://aclanthology.org/2021.iwslt-1.pdf#page=124>

RELATED WORKS

- End-to-End Speech Translation with Pre-trained Models and Adapters: UPC at IWSLT 2021



- M-Adapter: Modality Adaptation for End-to-End Speech-to-Text Translation



RELATED WORKS

- SALM: Speech-augmented Language Model with In-context Learning for Speech Recognition and Translation



<https://ieeexplore.ieee.org/abstract/document/10447553/>

RELATED WORKS

- SALM: Speech-augmented Language Model with In-context Learning for Speech Recognition and Translation



- BLSP: Bootstrapping Language-Speech Pre-training via Behavior Alignment of Continuation Writing



RELATED WORKS

- SALM: Speech-augmented Language Model with In-context Learning for Speech Recognition and Translation



- BLSP: Bootstrapping Language-Speech Pre-training via Behavior Alignment of Continuation Writing

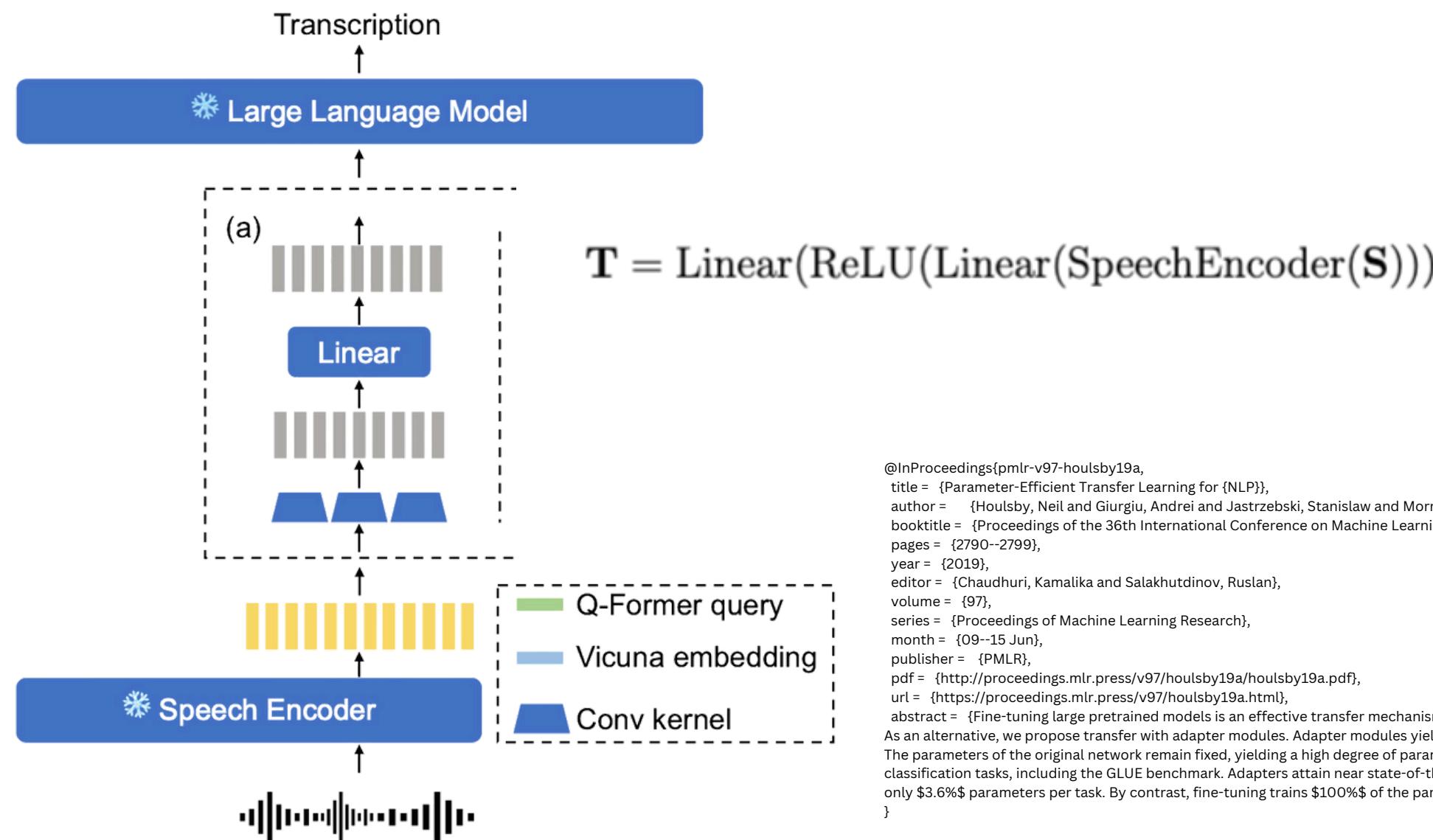


- SLM: Bridge the thin gap between Speech and Text foundation models



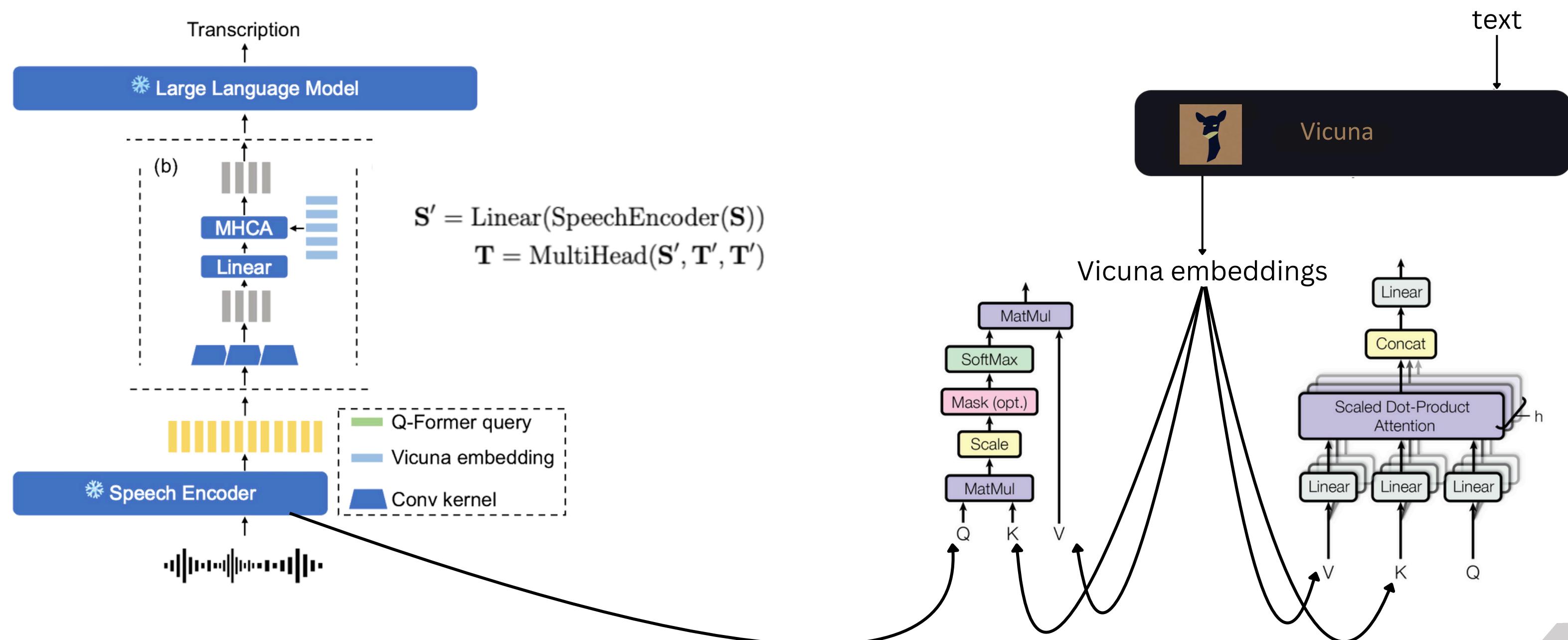
RELATED WORKS

- Connecting speech encoder and large language model for ASR



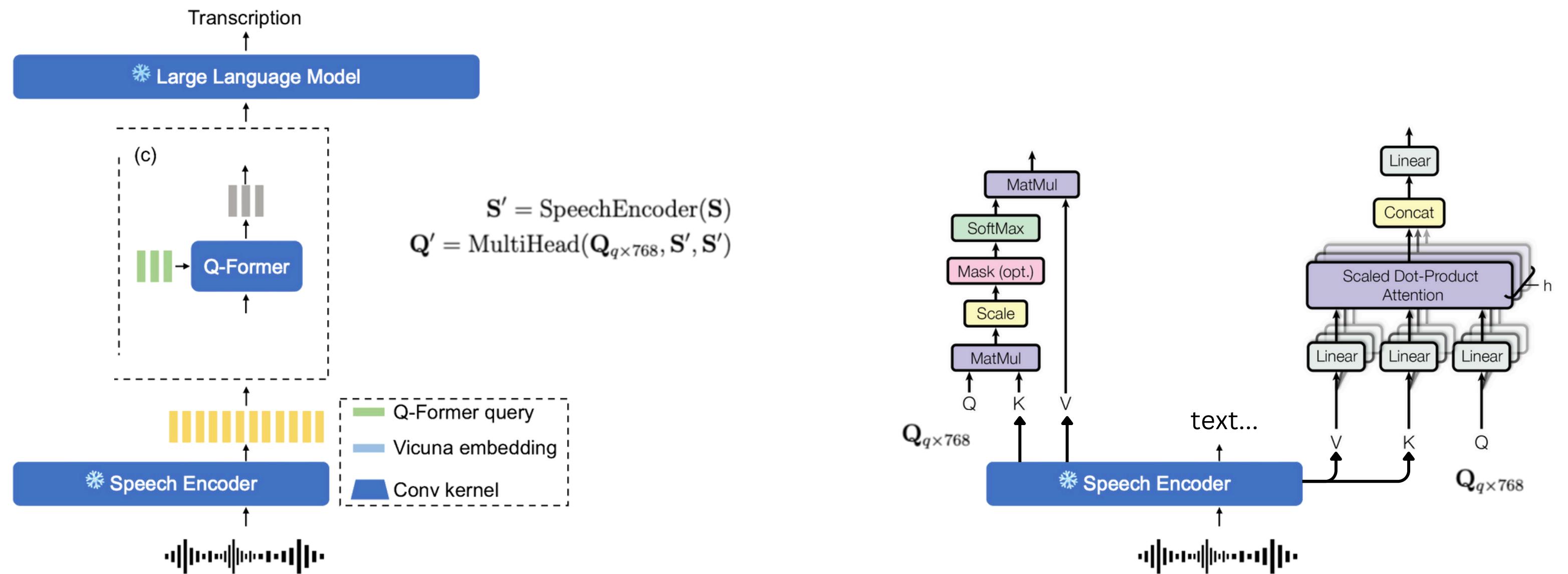
RELATED WORKS

- Connecting speech encoder and large language model for ASR



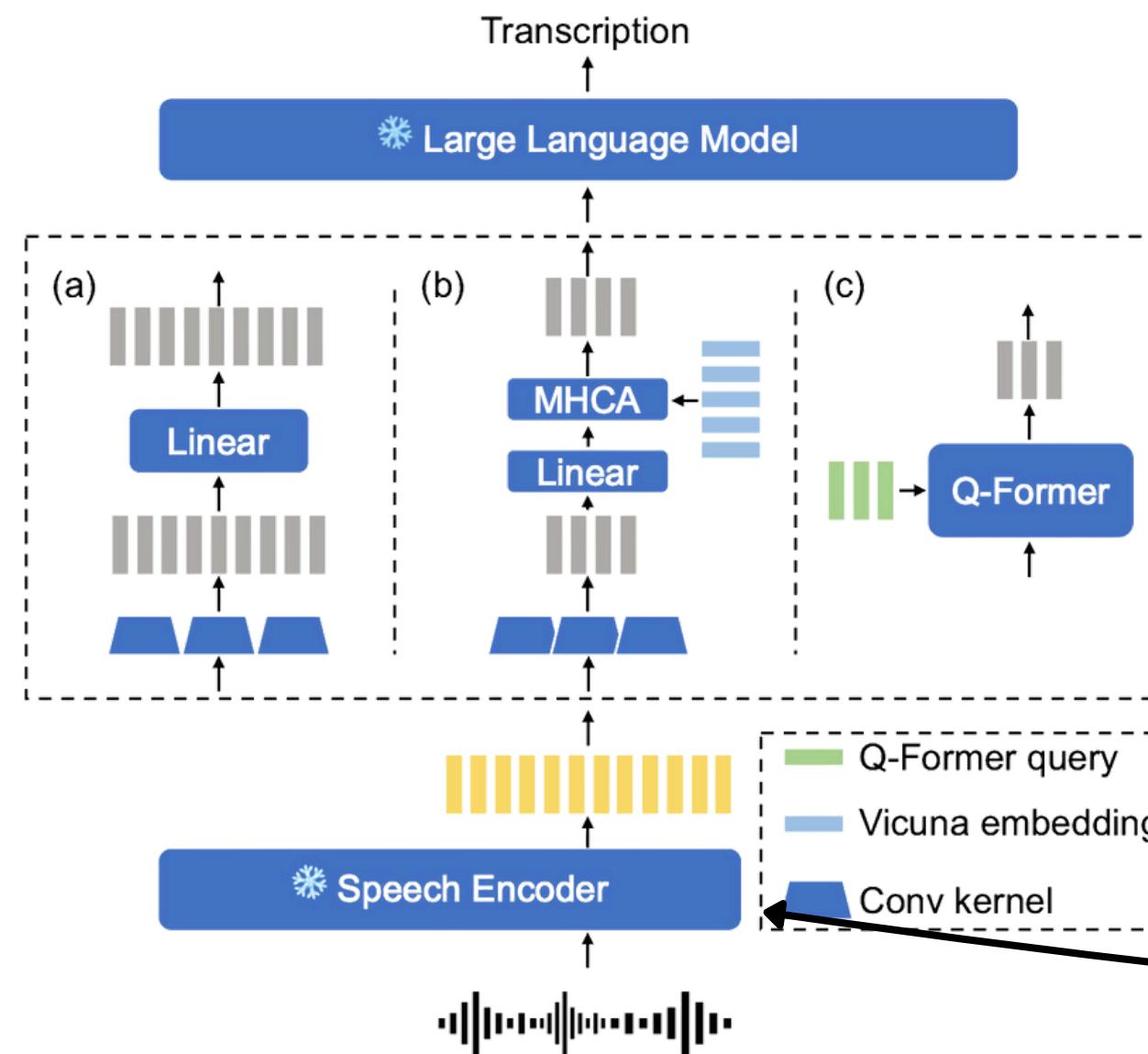
RELATED WORKS

- Connecting speech encoder and large language model for ASR



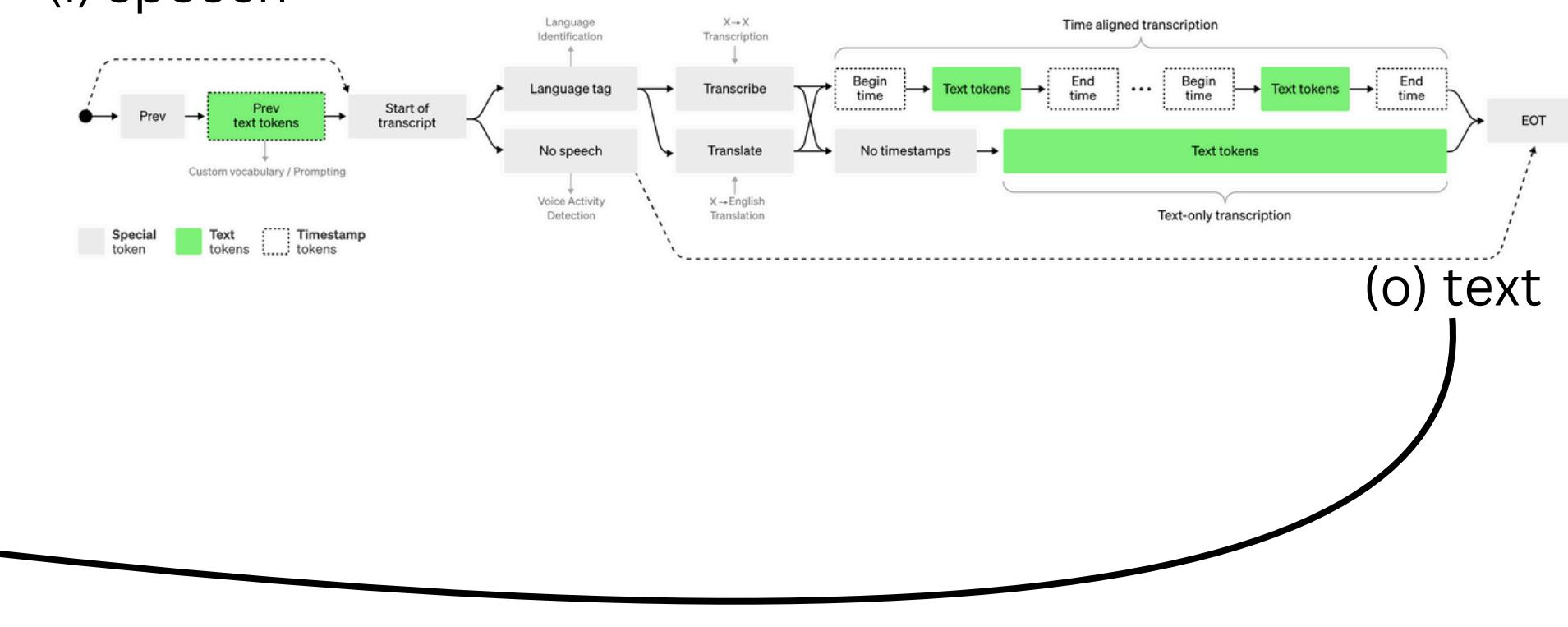
RELATED WORKS

- Connecting speech encoder and large language model for ASR



whisper speech encoder

(i) speech

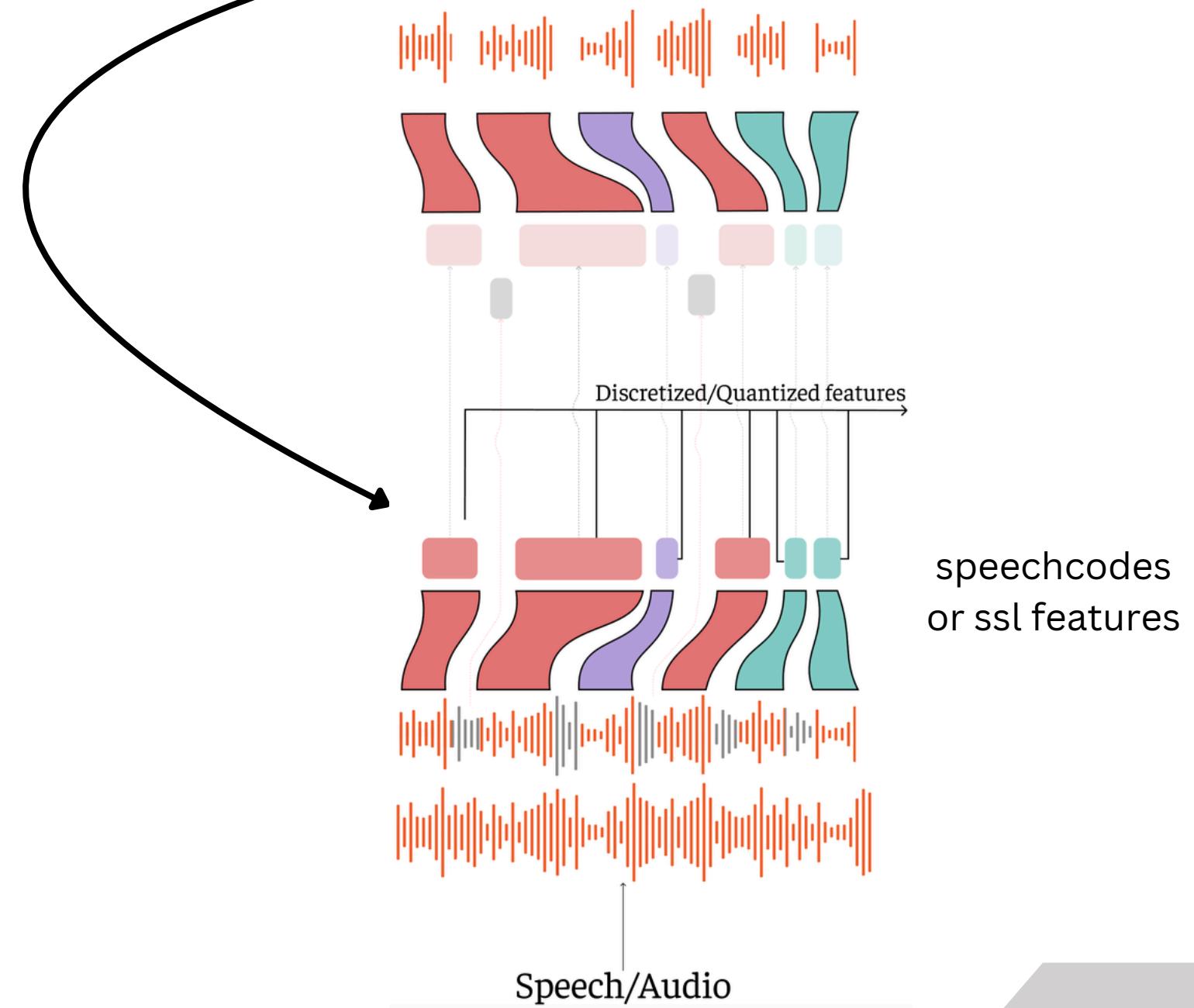


(o) text

BACKGROUND

- we'll try parameter efficient way by freezing LLM and speech encoder while training only modality connector/adapter

...high level idea



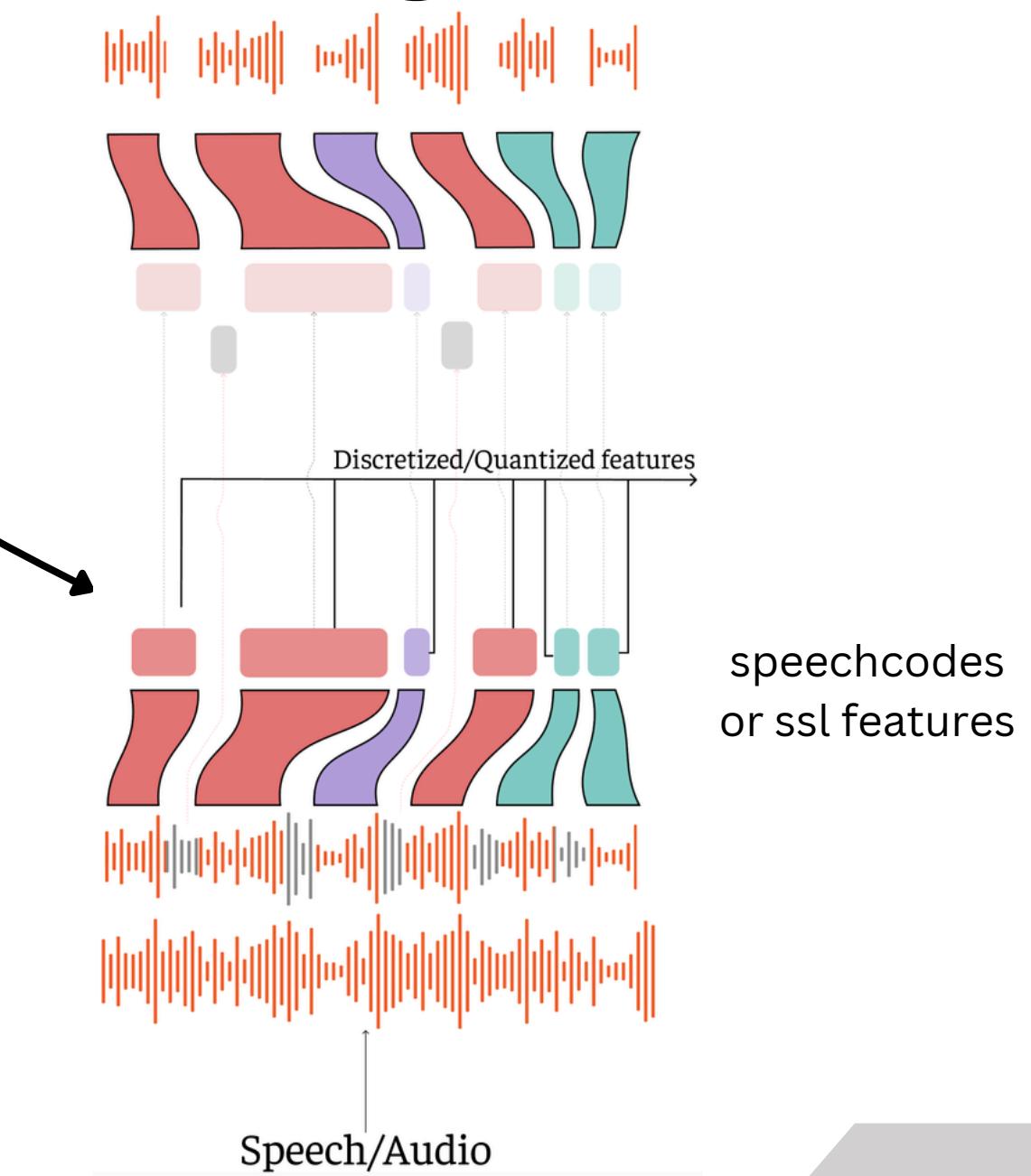
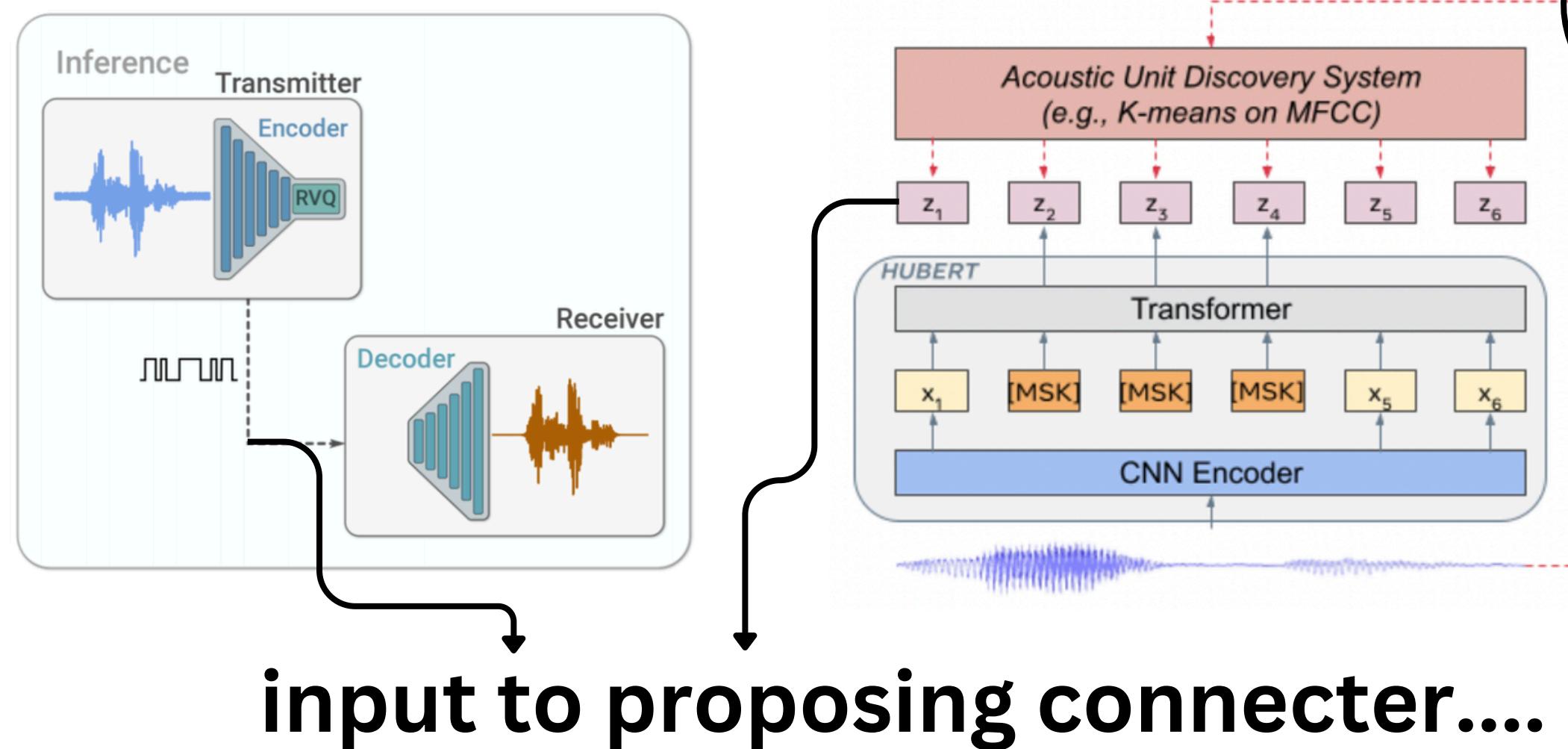
* ssl - self-supervised learning

BACKGROUND

- we'll try parameter efficient way by freezing LLM and speech encoder while training only modality connector/adapter

...high level idea

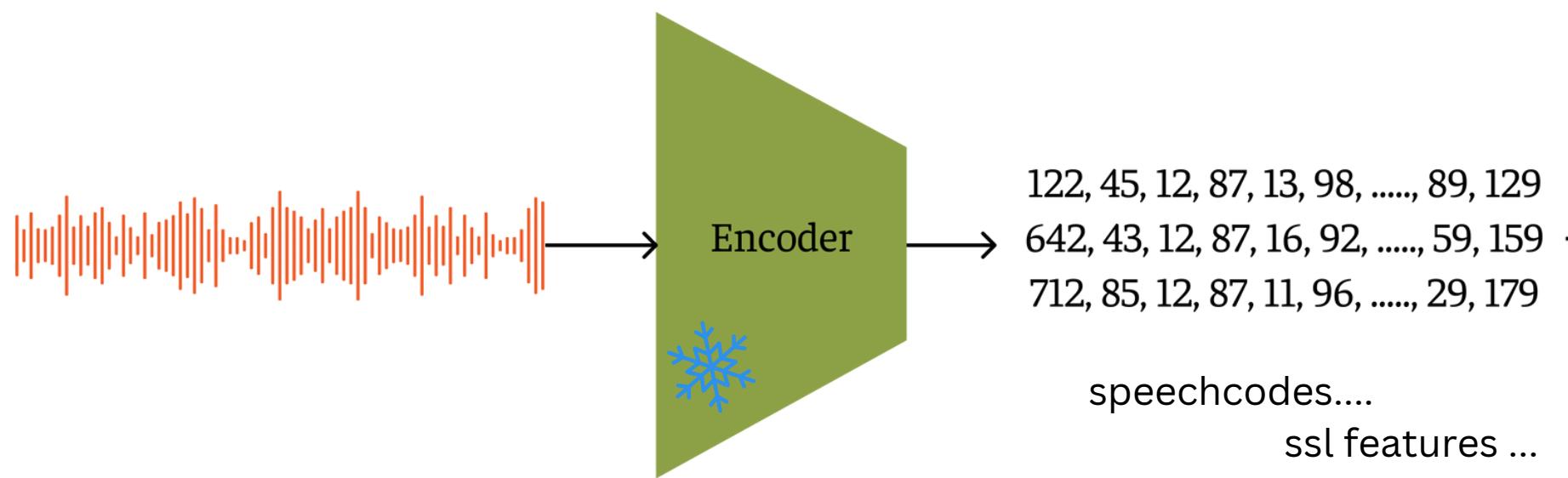
...detailed



* ssl - self-supervised learning

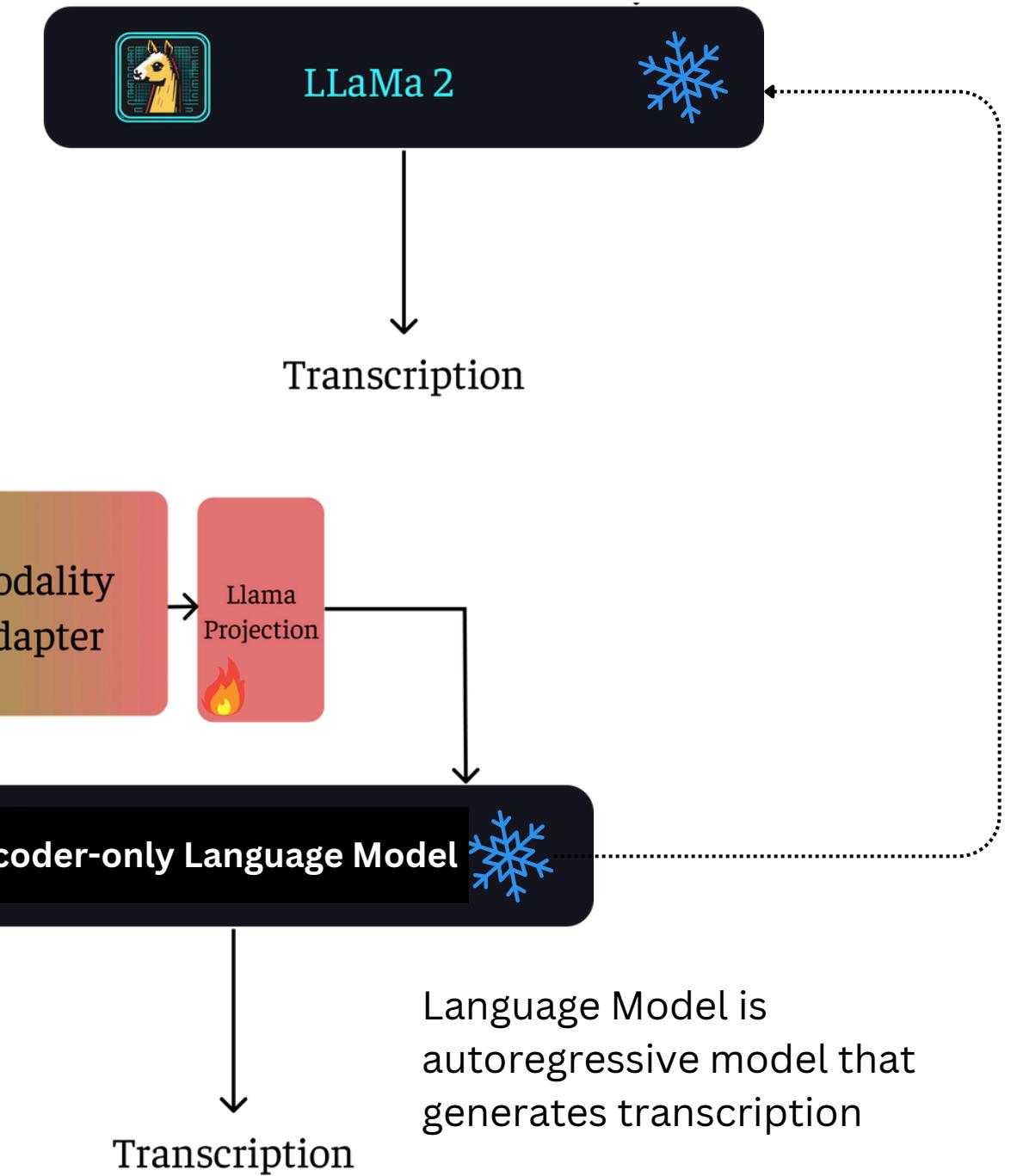
PROPOSED METHODS

modality adapter that aligns speech to te

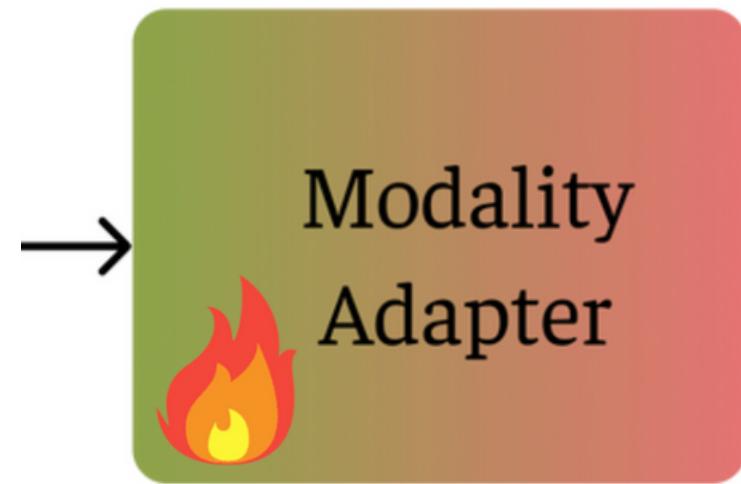


Speech Encoder --
Neural Encoderc (SpeechTokenizer)
<https://arxiv.org/abs/2308.16692>

Speech Encoder --
HuBERT-base, HuBERT-large, HuBERT-large-ft
<https://arxiv.org/abs/2106.07447>

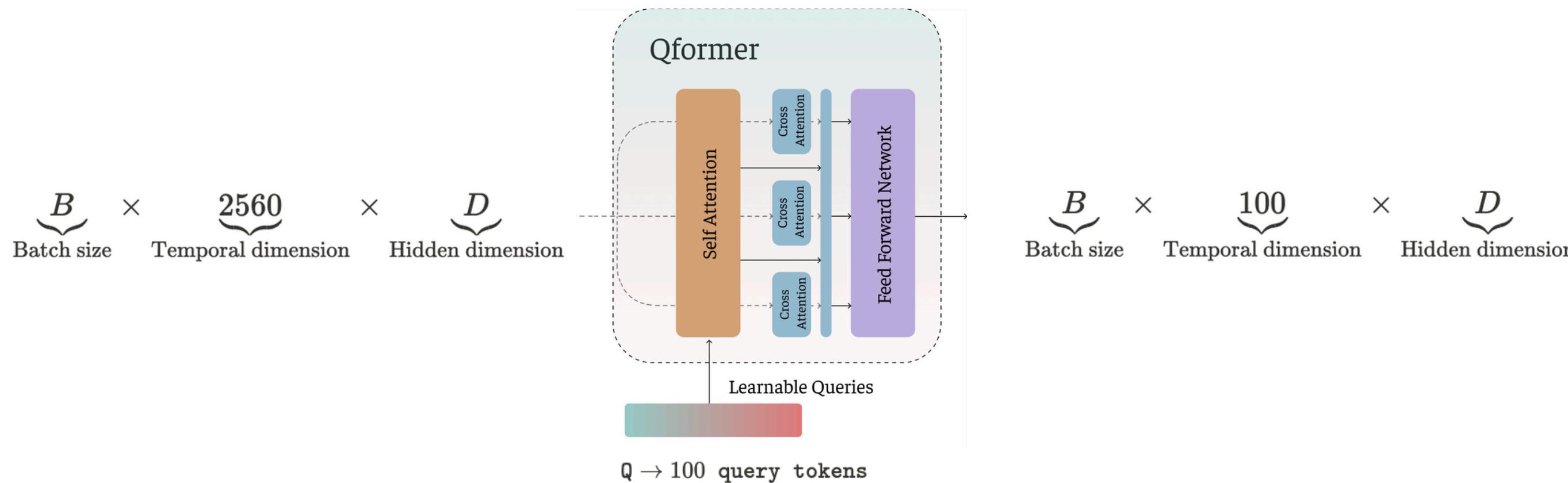


CURRENT WORKS

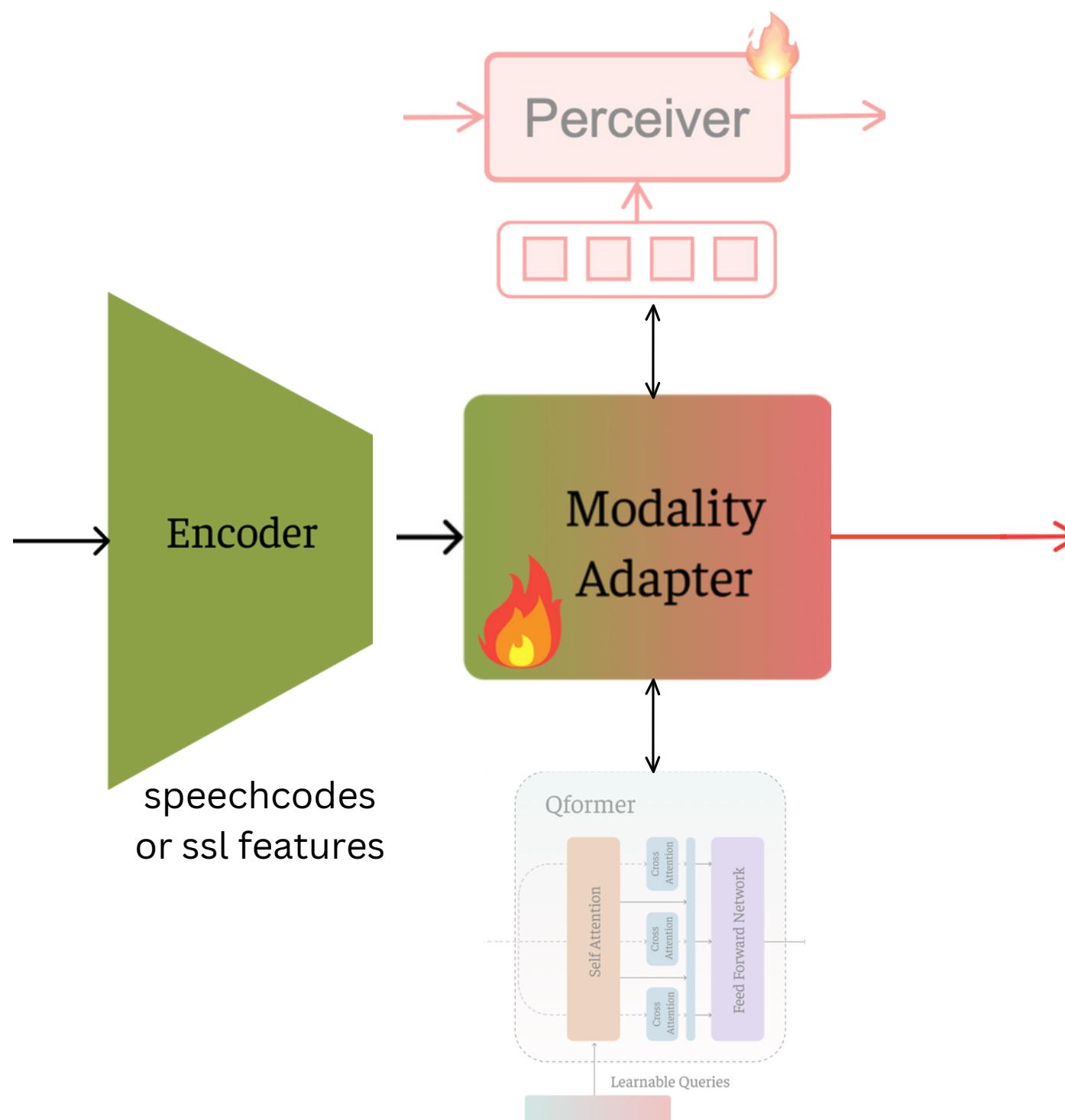


BLIP-2: Bootstrapping Language-Image
Pre-training with Frozen Image Encoders
and Large Language Models
<https://arxiv.org/abs/2301.12597>

QFormer - Querying Transformer

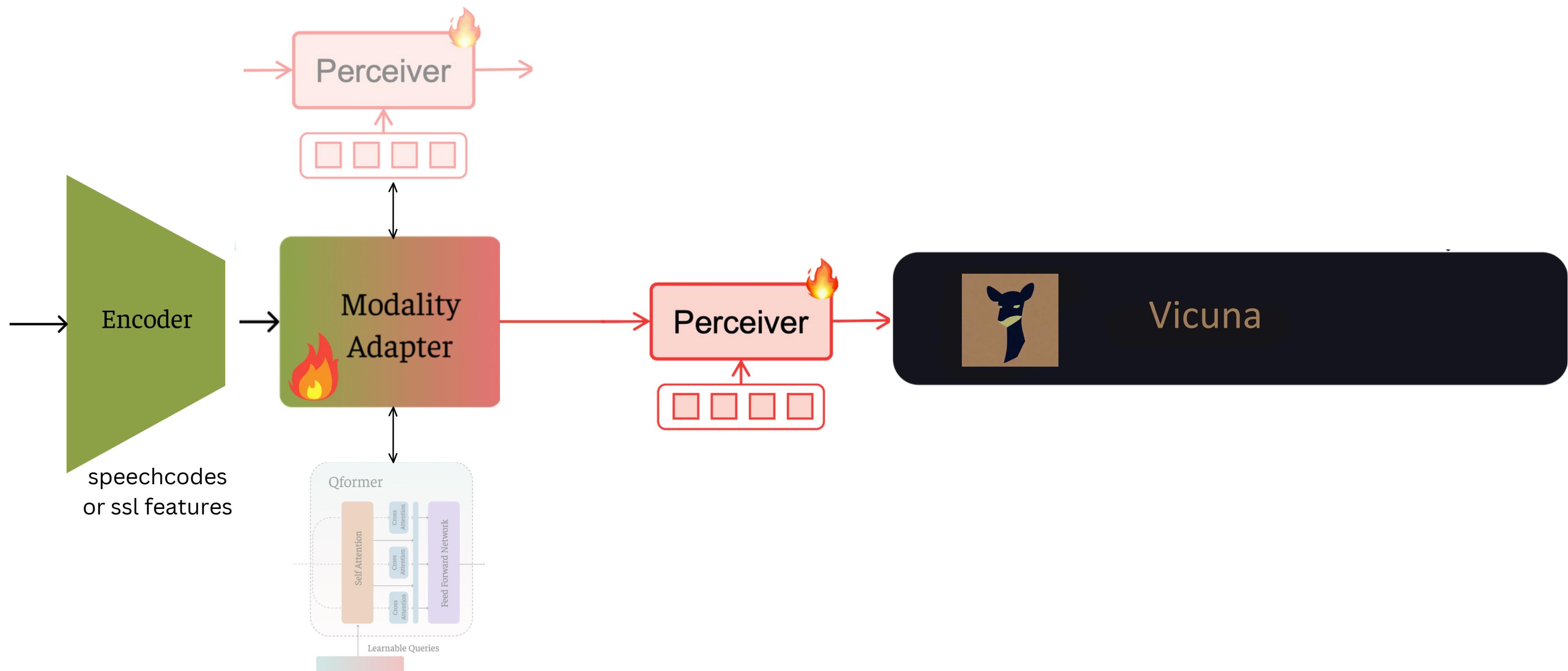


PROPOSED METHOD



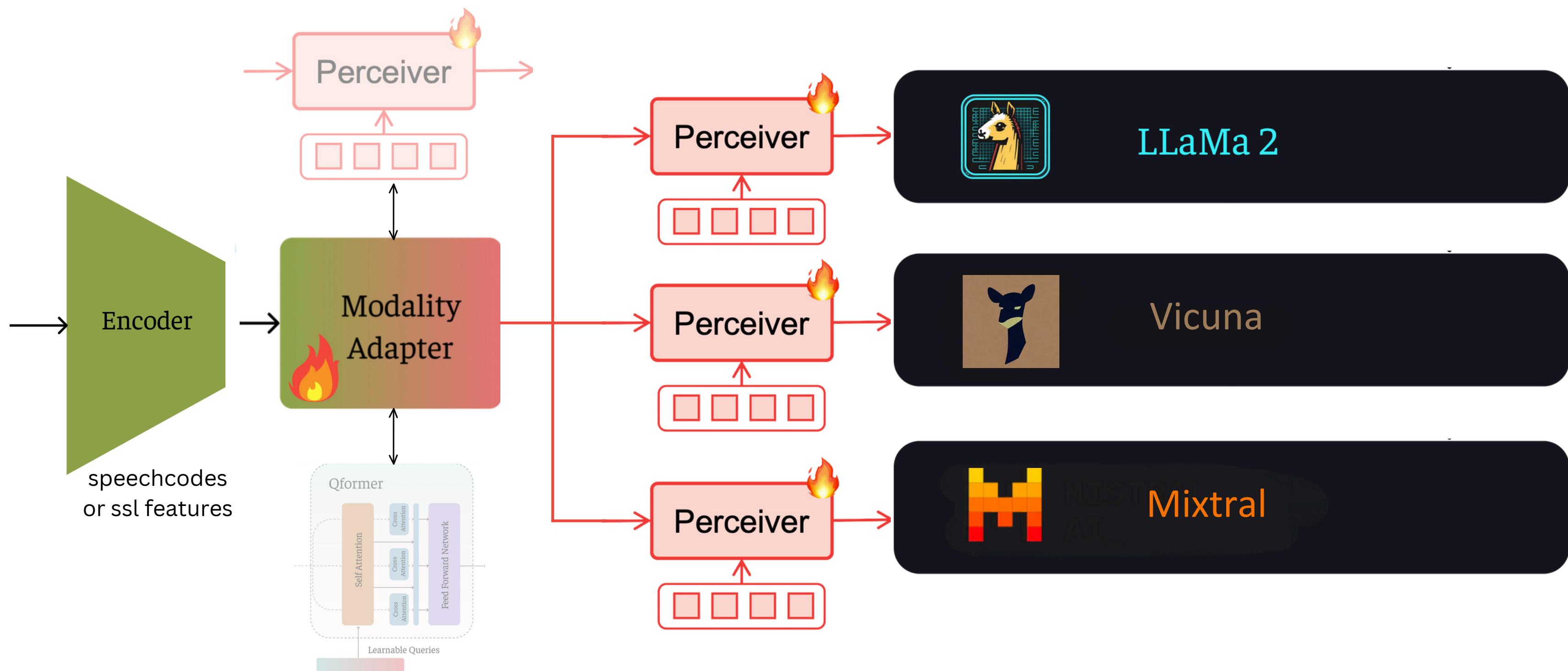
* Perceiver is a Modality adapter, predefined set of query tokens serve as adapter tokens

PROPOSED METHOD



* Perceiver is a Modality adapter, predefined set of query tokens serve as adapter tokens

PROPOSED METHOD



* Perceiver is a Modality adapter, predefined set of query tokens serve as adapter tokens

DATA PREPROCESSING

Table 4.1: Table Example: Comparative analysis of target transcriptions versus predictions after processing through a neural EnCodec.

Nº	Target Transcription	Predicted Transcription
(1)	and now the squires despatched to make the prize came up bringing with them two gentlemen on horseback two pilgrims on foot and a coach full of women with some six servants on foot and on horseback in attendance on them and a couple of muleteers whom the gentlemen had with them	and now the squires dispatched to make the prize came up, bringing with them two gentlemen on horseback two pilgrims on foot, and a coachful of women with some six servants on foot and on horseback in attendance on them and a coachful of women with some six servants on foot and on horseback in attendance on them and a coachful of women with some six servants on foot and on horseback in attendance on them and a coachful
(2)	otto winked at me	i don't win to me
(3)	sonya smiled no	saw your smiles no
(4)	daphne exclaimed tenderly what more is needed	to have me explain tenderly what morrisley did
(5)	they look too funny for anything the opossums are the only marsupials in this country continued old mother nature (non-English symbols)
(6)	our prostrate brother homo ventrambulans razor	a prostrate brother a moe then trampolines razor
(7)	pog n	pog now
(8)	repugnant alike to reason and conscience resigned to growing infirmities resist a common adversary	they don't end up alike, they're reasoned and conscious. redesigned, they're cruelly made in their medias. they're reasoned, they come in adversity
(9)	for a sound that would mean the near approach of danger hush sh sh came again as a gentle murmur from below and the something that moved and breathed in the darkness seemed to draw nearer to yvonne	for a sound that would main the nearer put to danger hush shh shh came again as a gentle mammoth for below and the something that moved and breathed in the darkness seemed to draw nearer to yvonne

Assess the reconstruction capabilities of the speech encoder.

Whisper Large, one of the most precise commercial options for speech recognition, resulted in 22.5%

Further, text processing on predicted transcription that involves punctuation removal, normalizing letters reduced WER to 10.5% and CER to 5.6%

DATA PREPROCESSING

Table 4.1: Table Example: Comparative analysis of target transcriptions versus predictions after processing through a neural EnCodec.

Some predicted transcription examples may not closely resemble the actual transcription:

- (2), (3), and (8) examples

There was cases where ASR model Failed to recognize

- (5) example
....why?!

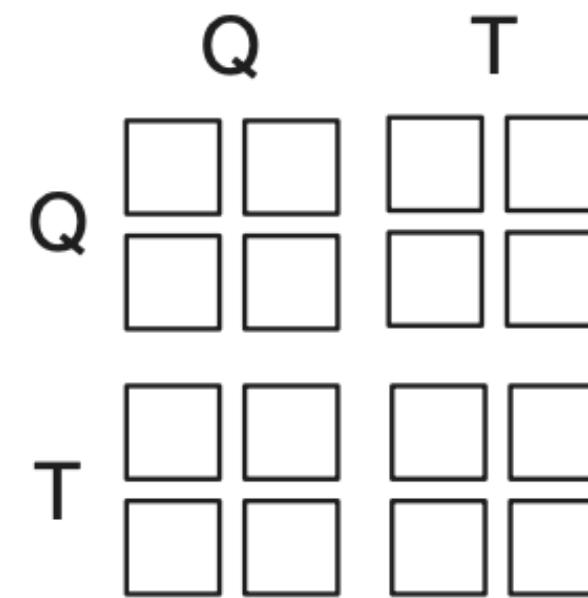
no paralinguistic features

- (7) example. However,
 - this was not always the case,
 - (9) example, which retained its paralinguistic elements.

TRAINING OBJECTIVE

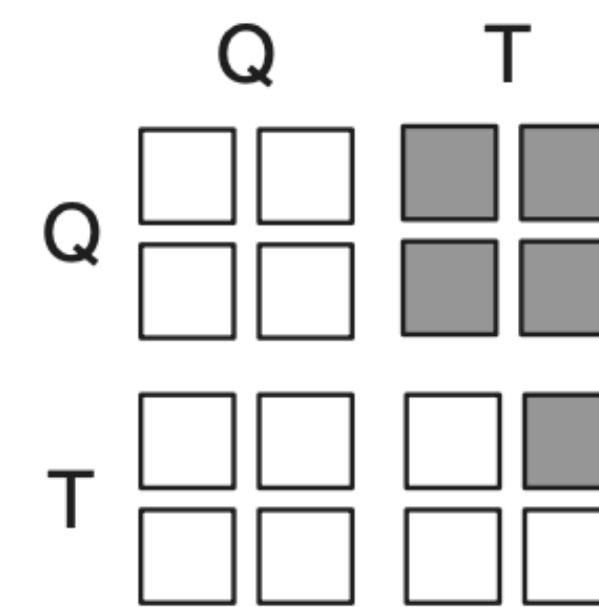
Pre-Train Speech-Text Qformer

Q: query token positions; **T:** text token positions.
■ masked □ unmasked



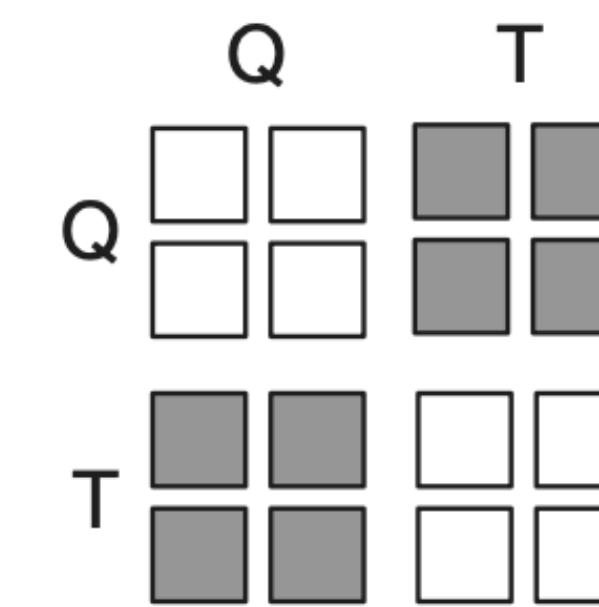
Bi-directional
Self-Attention Mask

Speech-Text
Matching



Multi-modal Causal
Self-Attention Mask

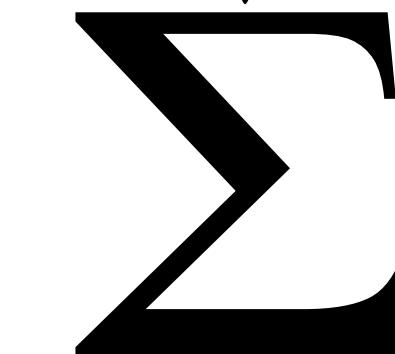
Speech-Text
Generation



Uni-modal
Self-Attention Mask

Speech-Text
Contrastive Learning

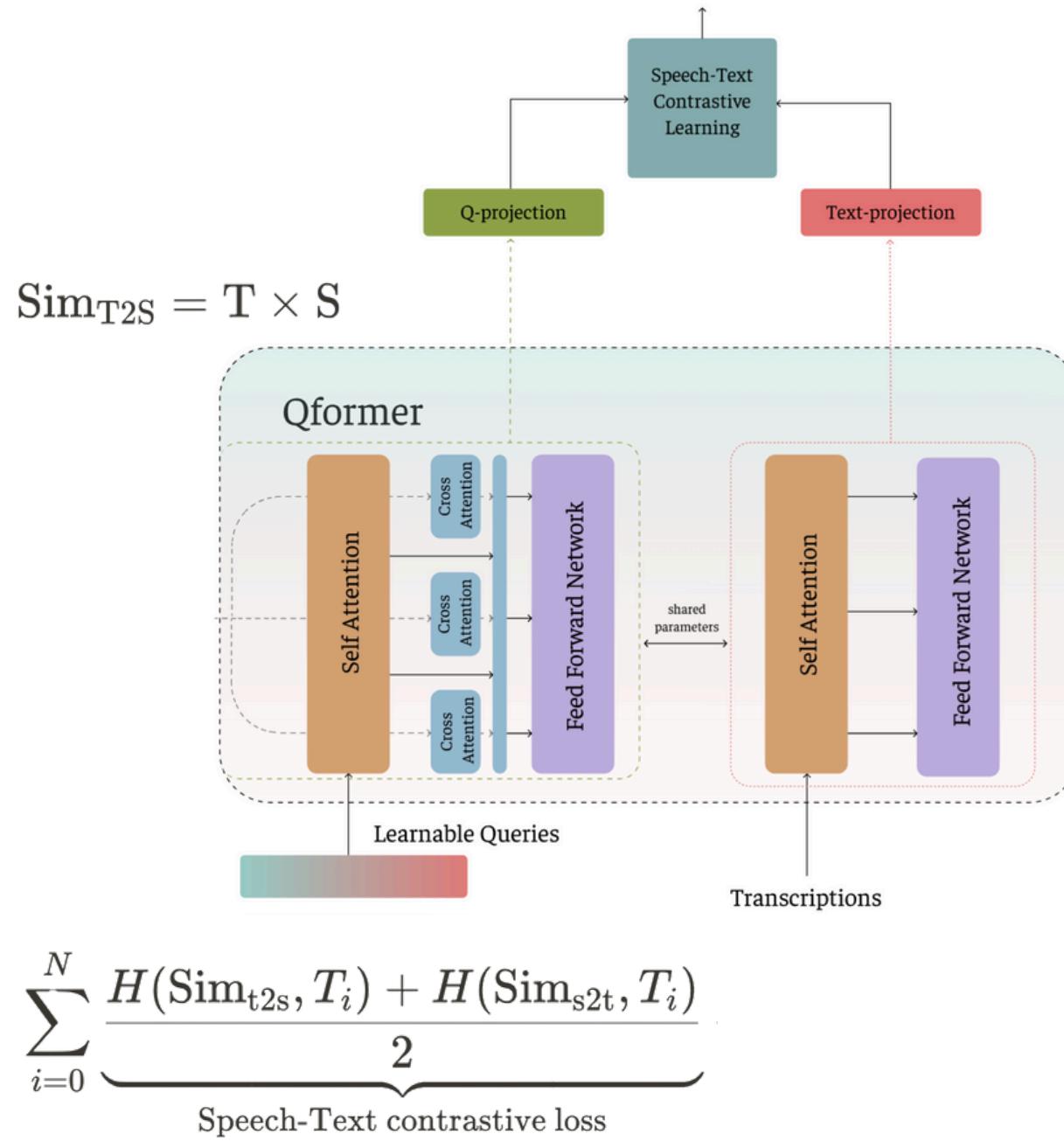
3 training losses: **STM**, **SGT**, **STCL**



STM - Speech Text Matching
SGT - Speech-Text Generation
STCL - Speech-Text Contrastive Learning

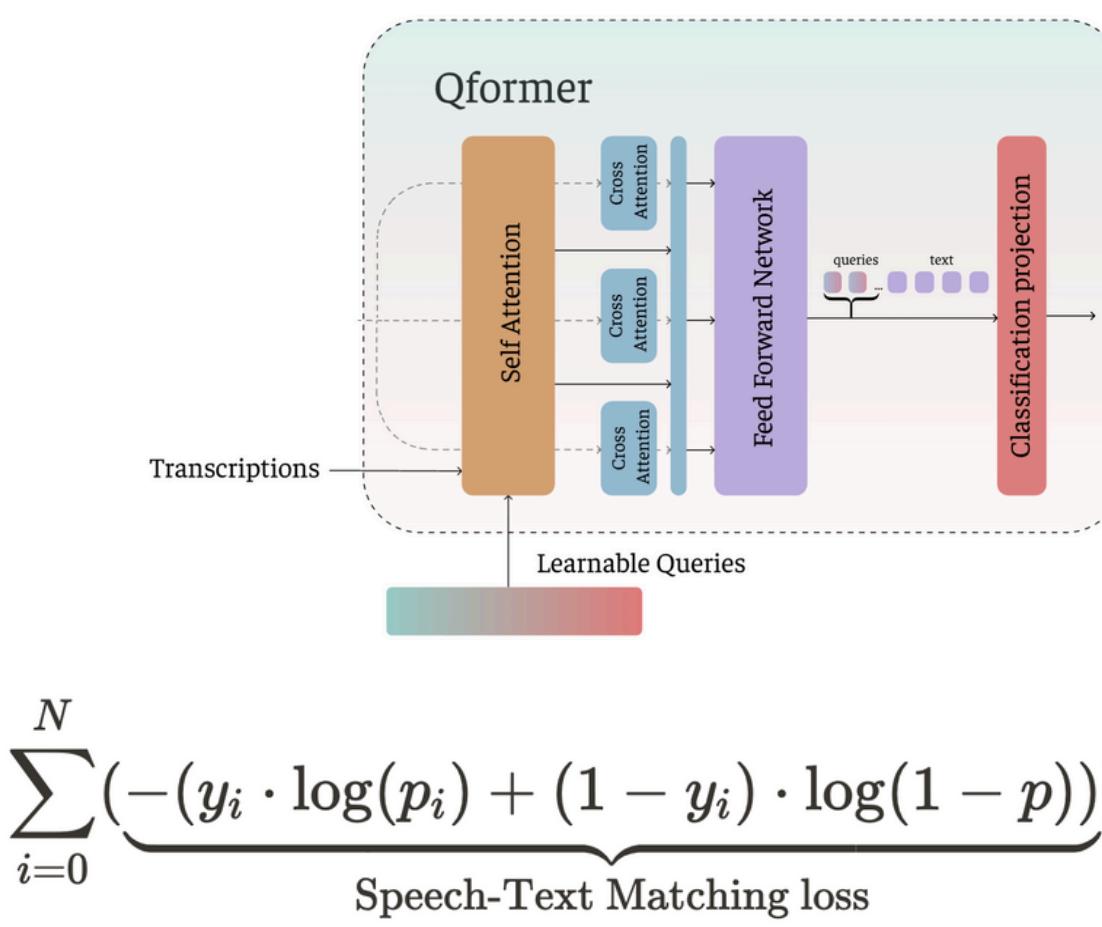
TRAINING OBJECTIVE

“Speech-Text Contrastive Training”



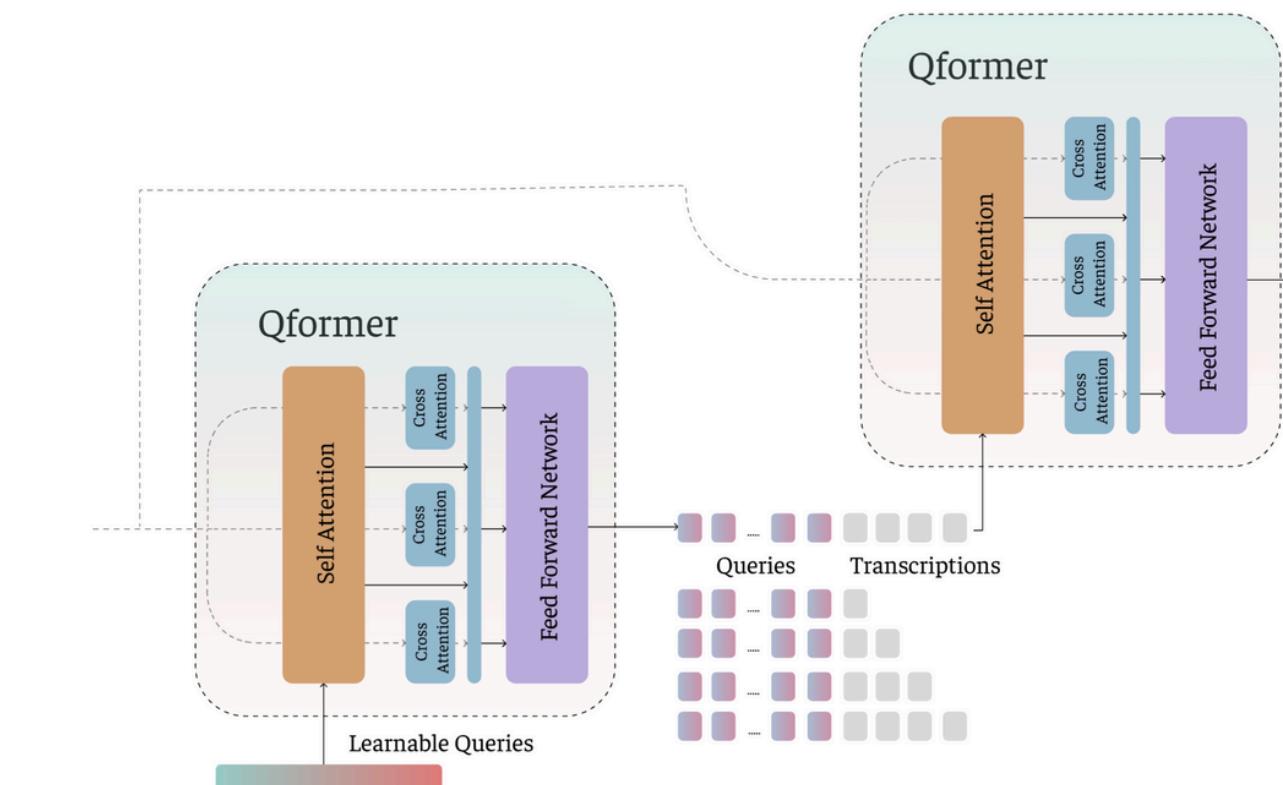
TRAINING OBJECTIVE

“Speech-Text Matching Training”



TRAINING OBJECTIVE

“Speech-Text Generation Training”

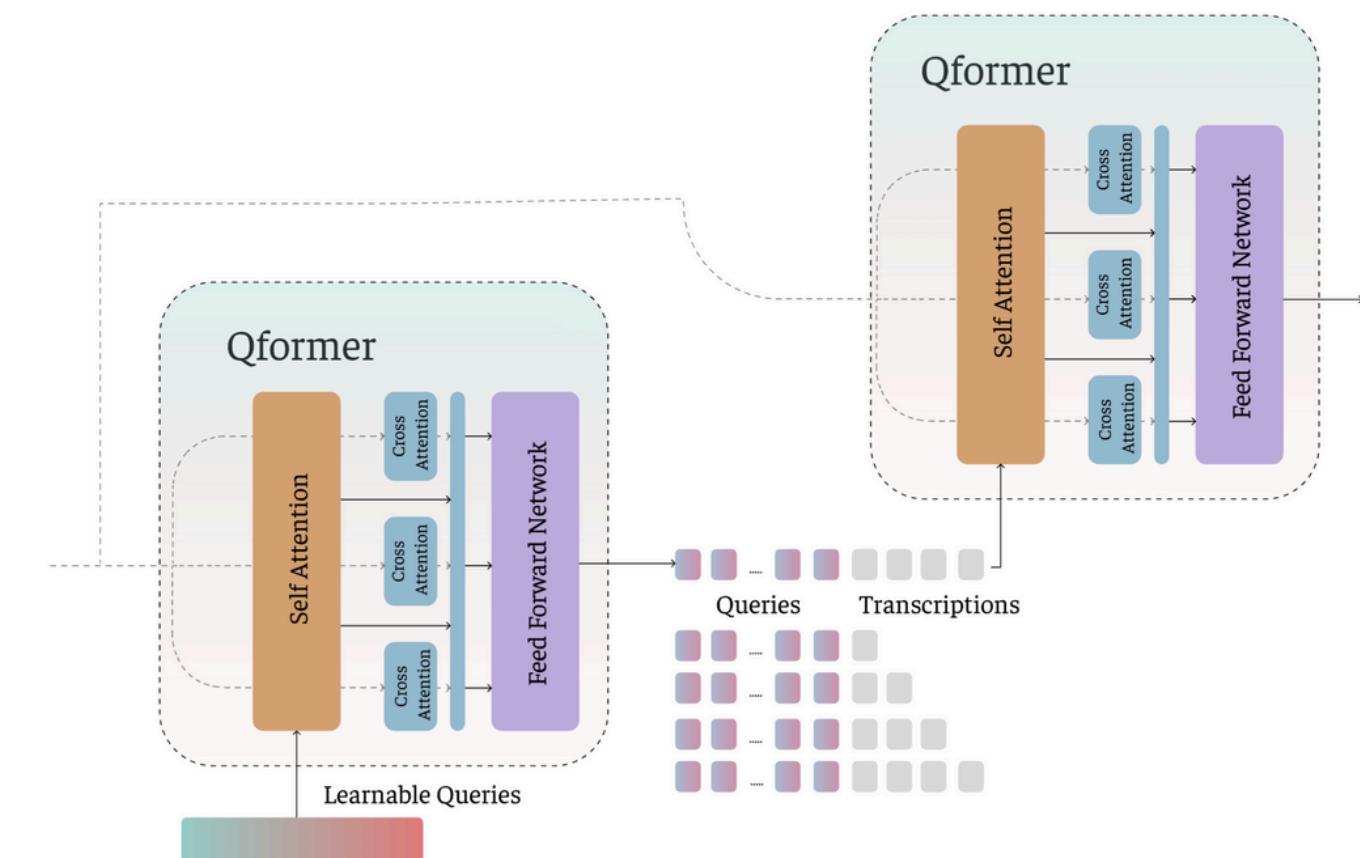
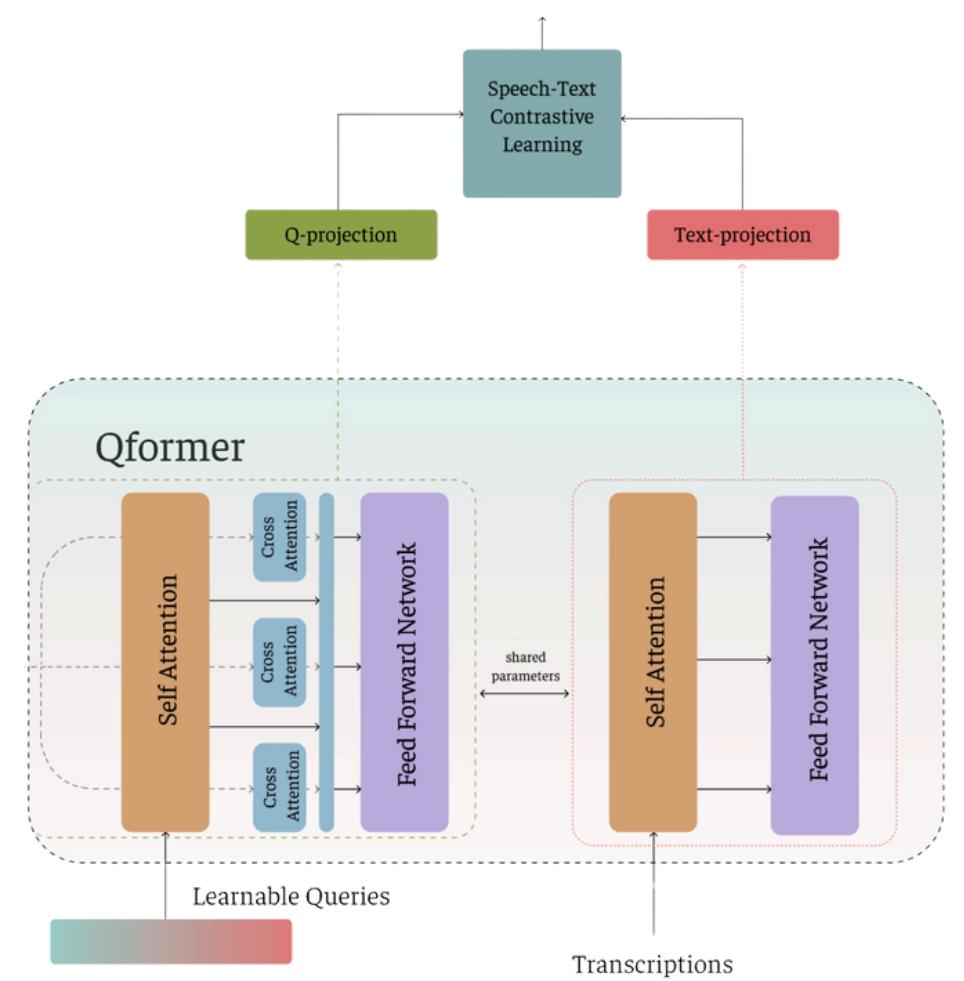


$$\sum_{i=1}^N \log (\mathbb{P}(u_i | u_{i-k}, \dots, u_{i-1}, \Theta))$$



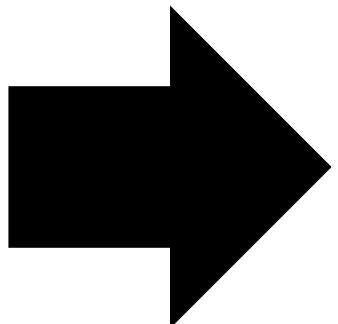
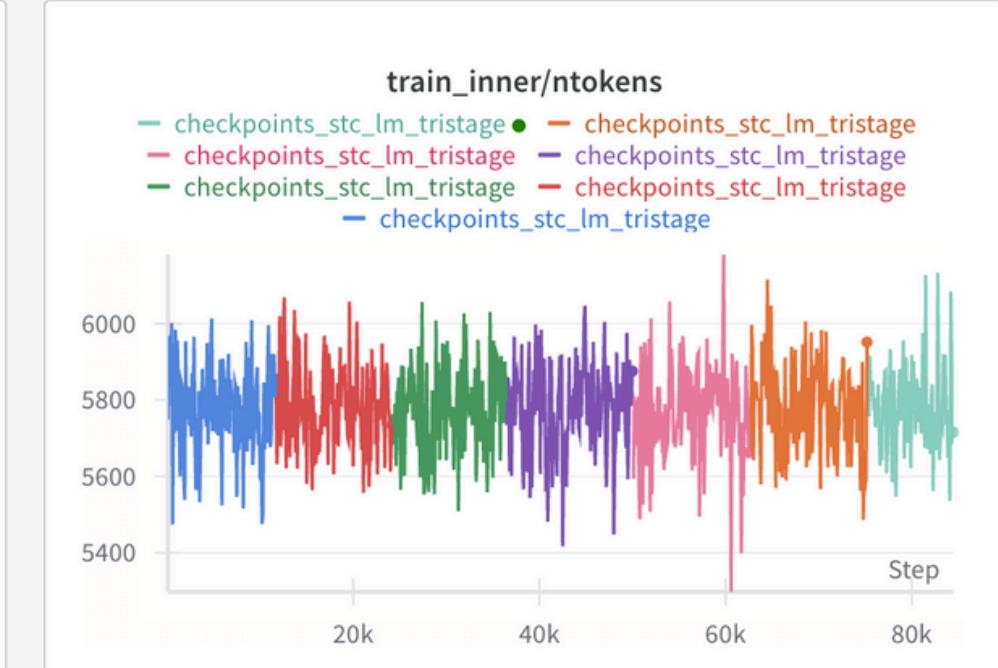
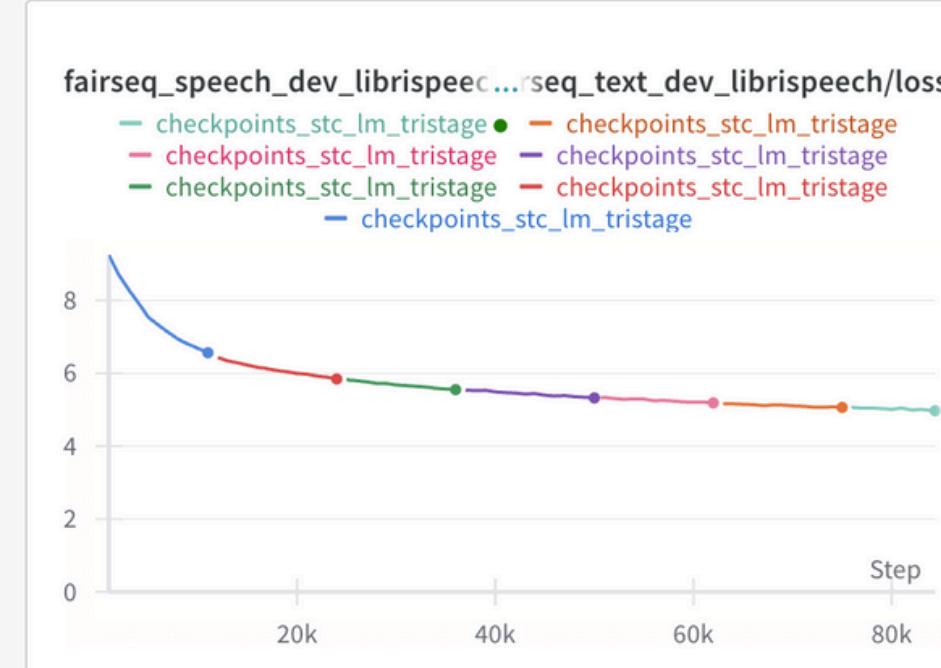
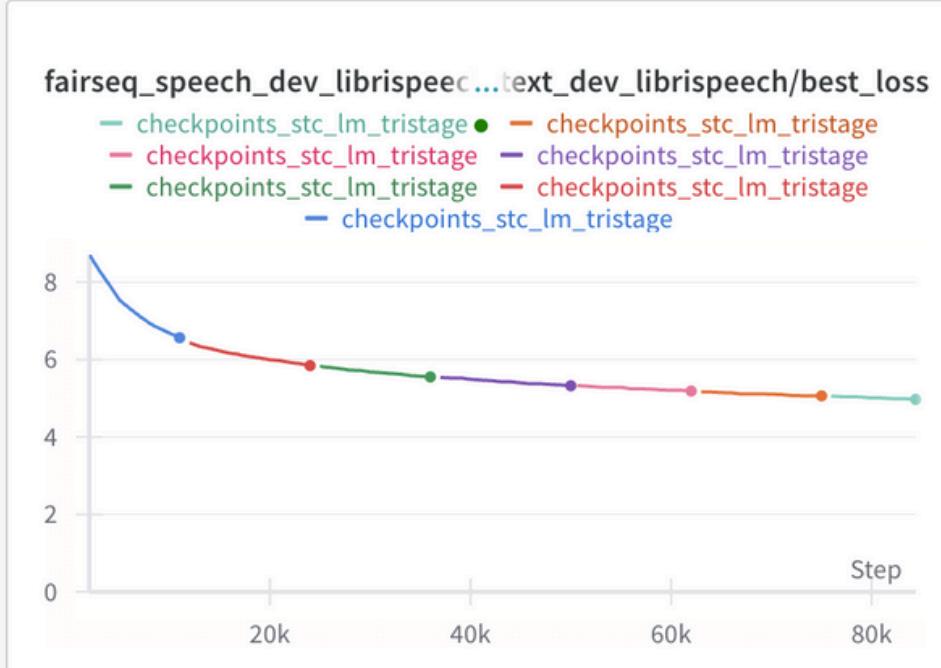
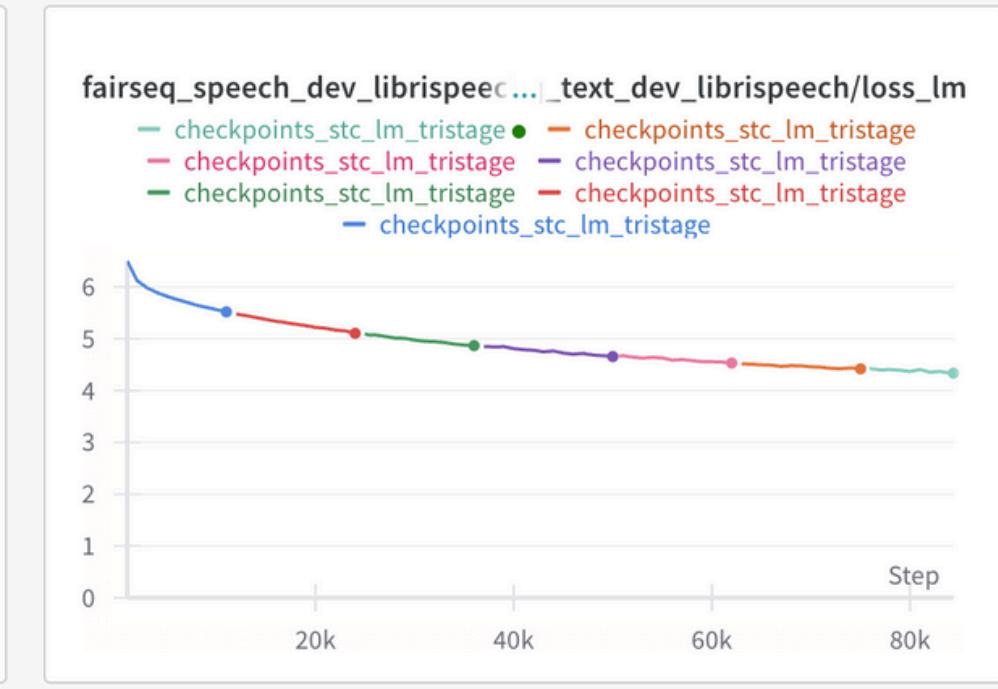
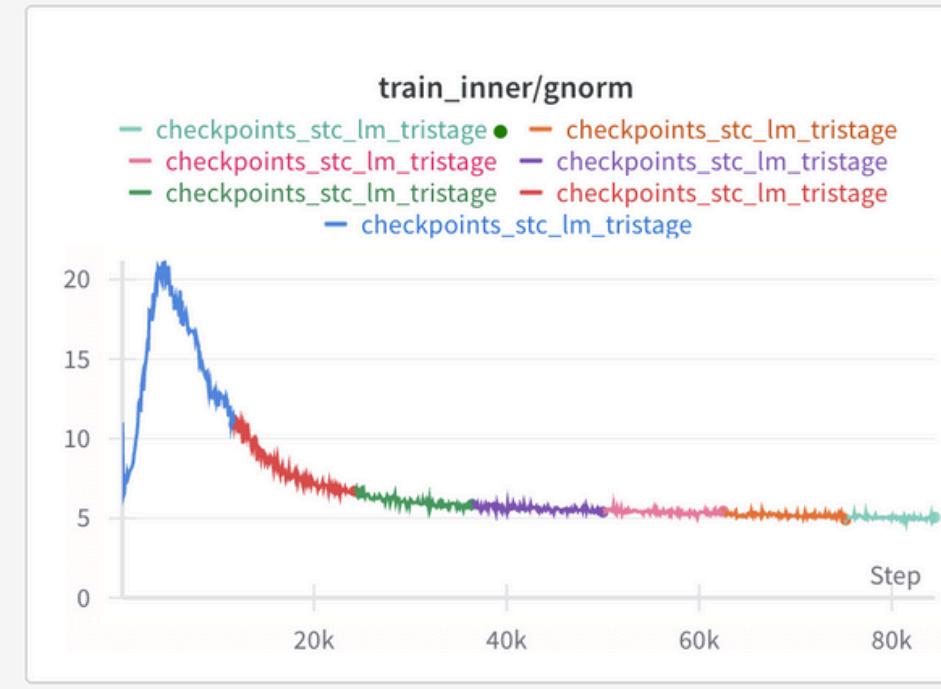
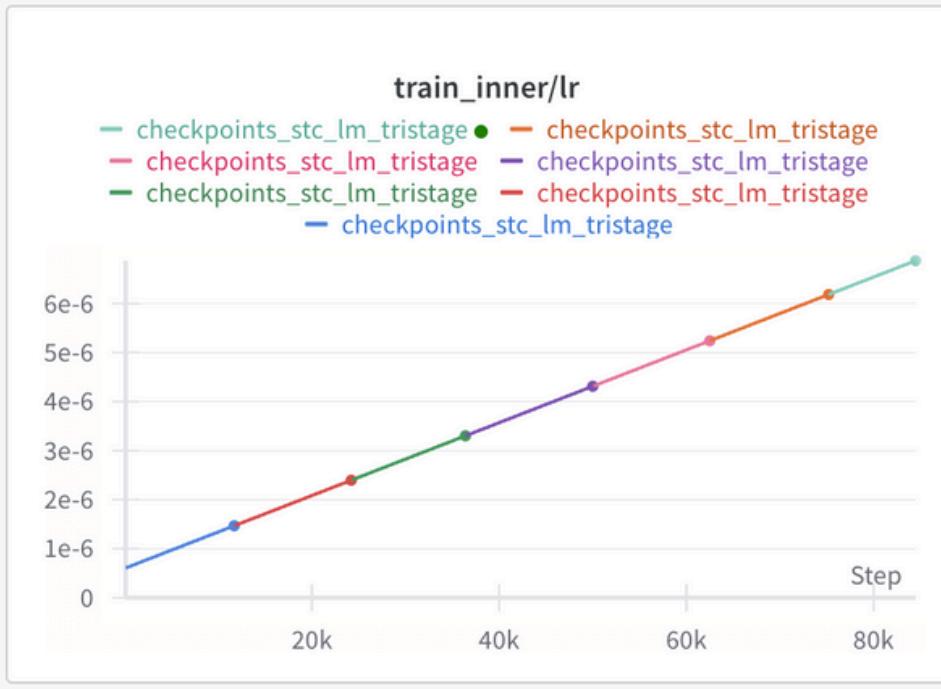
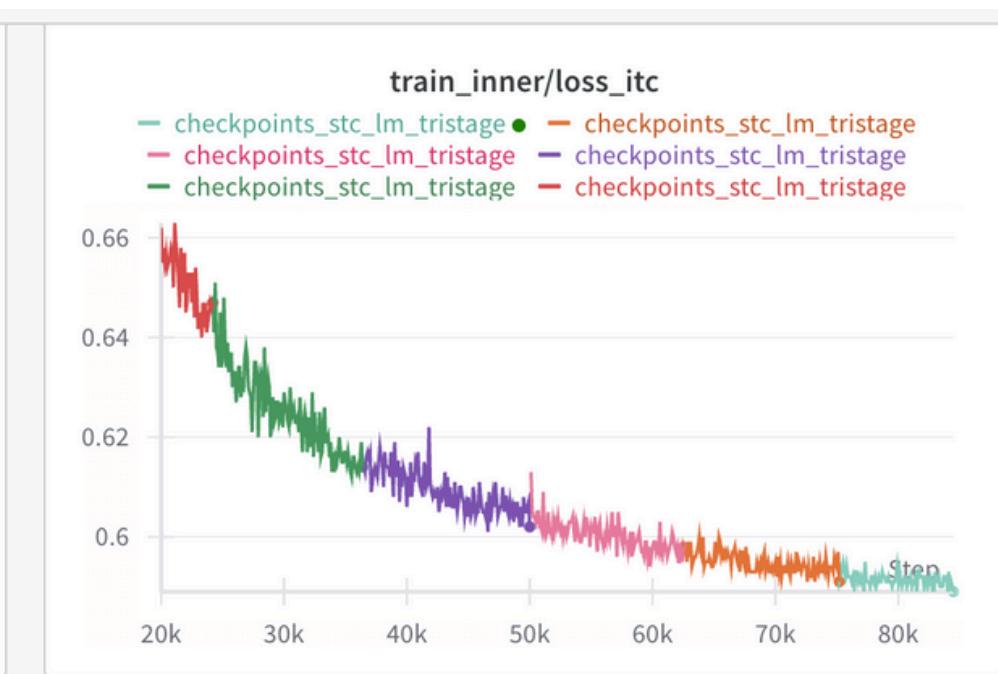
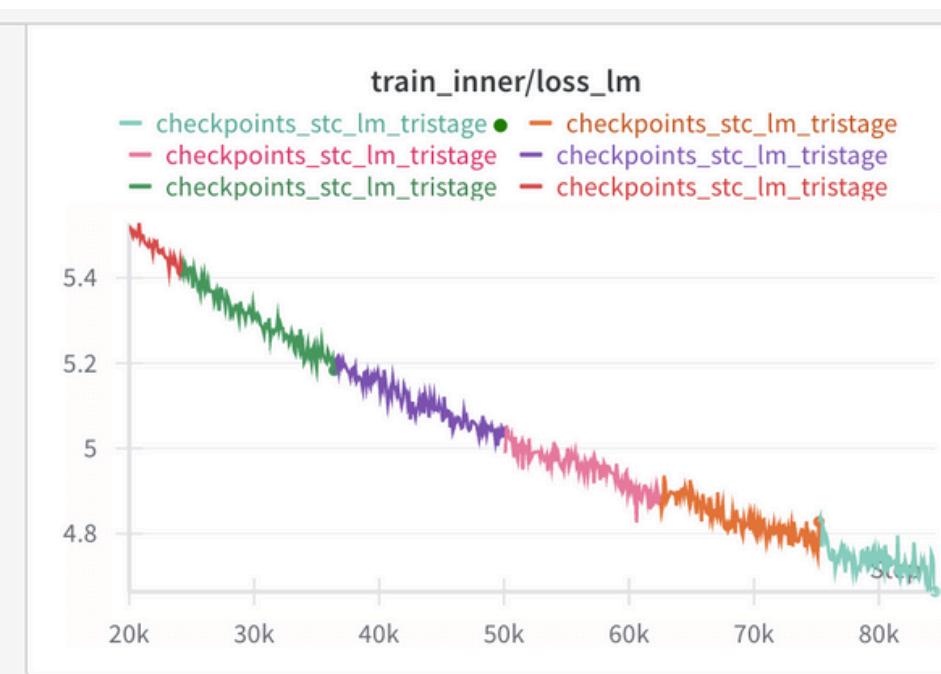
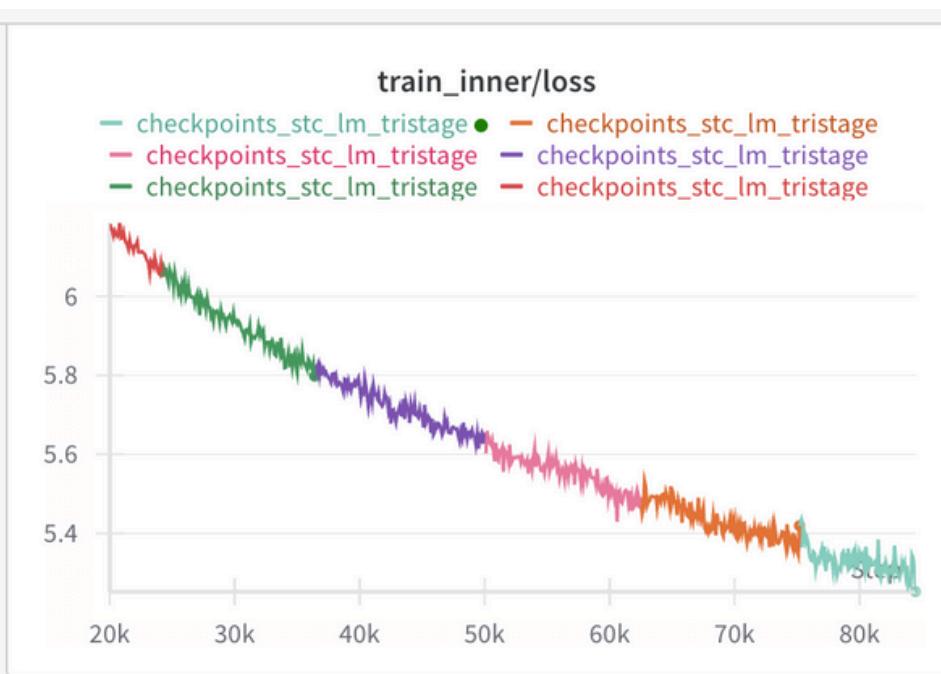
TRAINING OBJECTIVE

“Joint Speech-Text Contrasting & Speech-Text Generation Training”



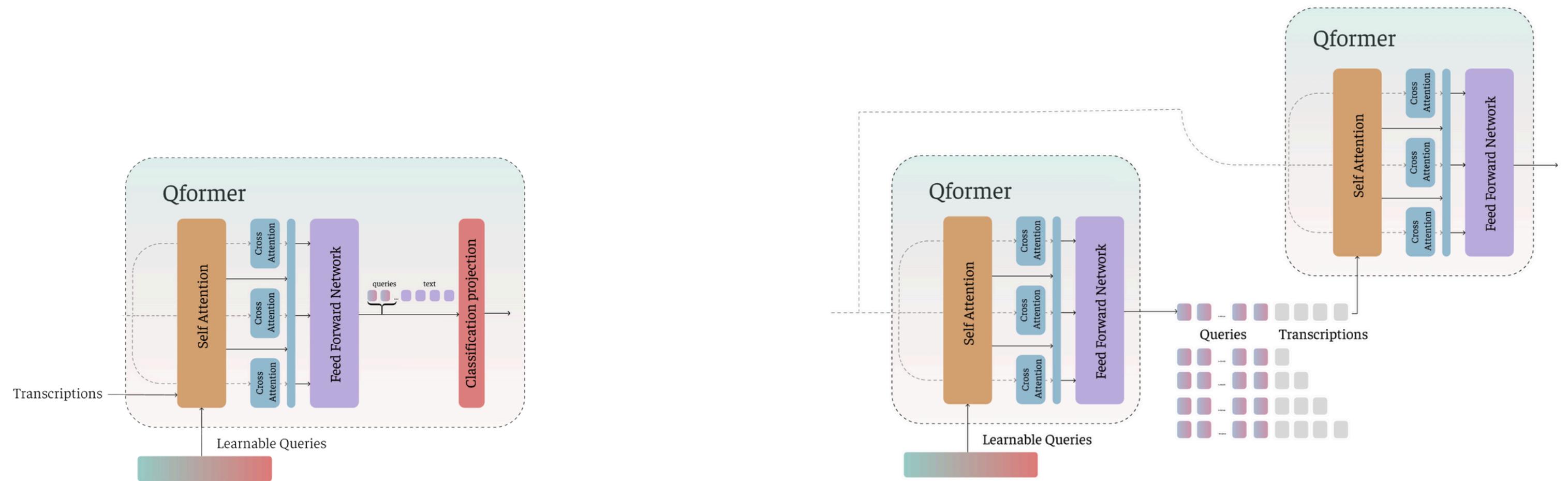
$$\sum_{i=0}^N \underbrace{\frac{H(\text{Sim}_{t2s}, T_i) + H(\text{Sim}_{s2t}, T_i)}{2}}_{\text{Speech-Text contrastive loss}}$$

$$+ \sum_{i=1}^N \log (\mathbb{P}(u_i | u_{i-k}, \dots, u_{i-1}, \Theta))$$



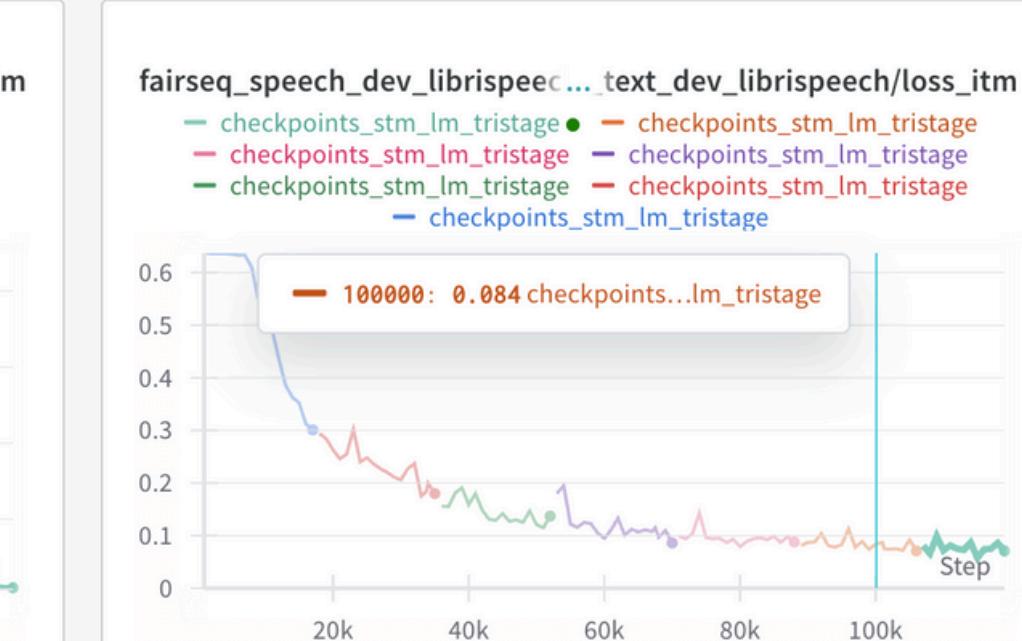
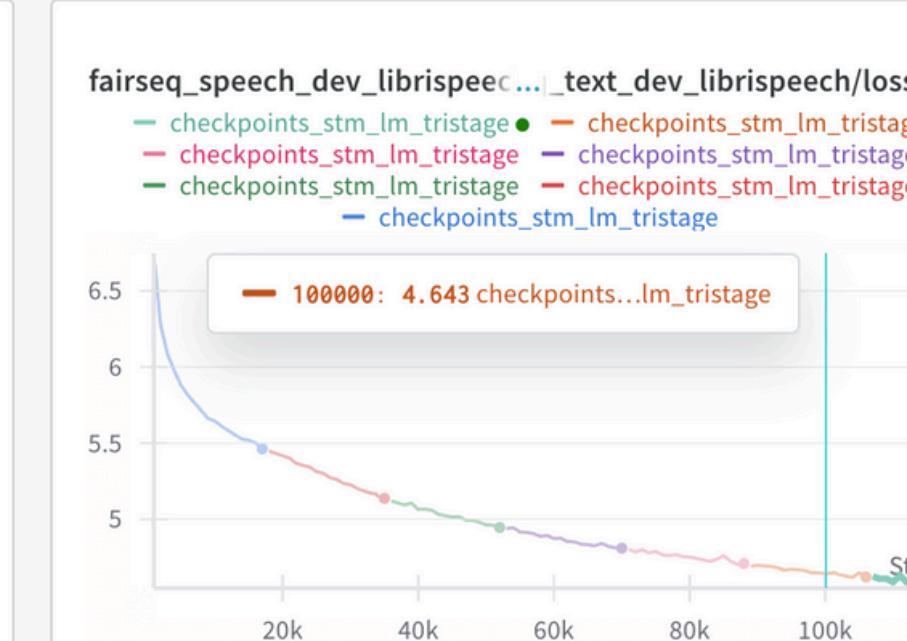
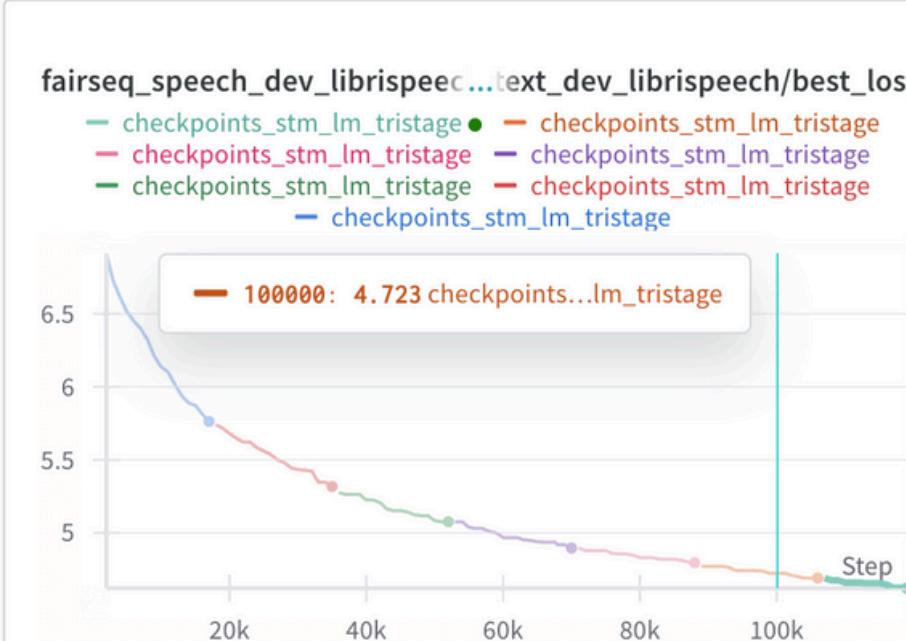
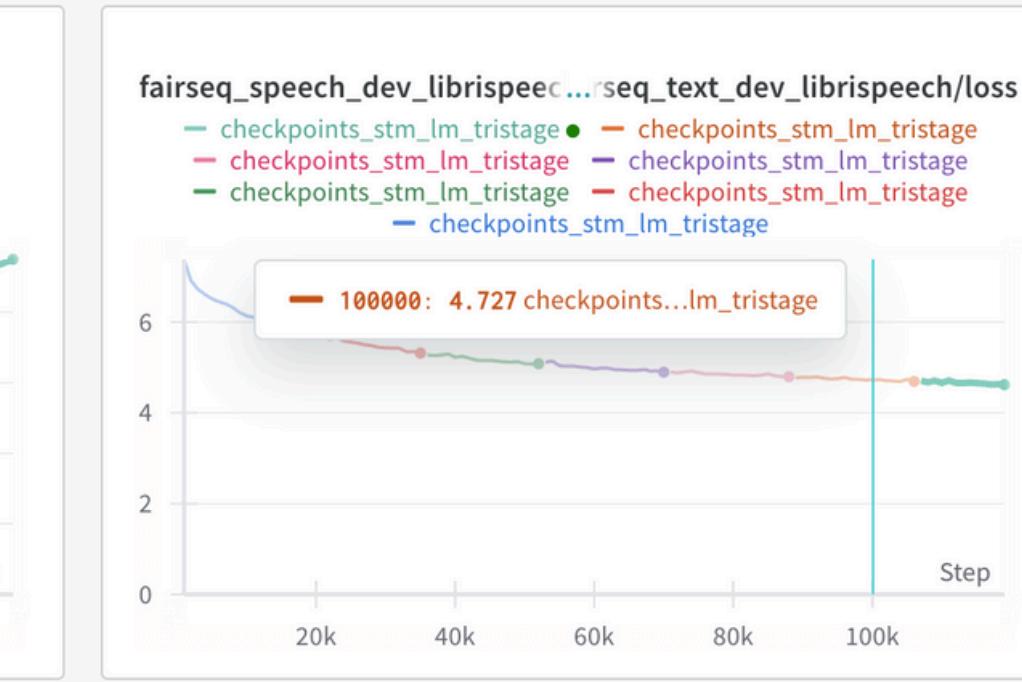
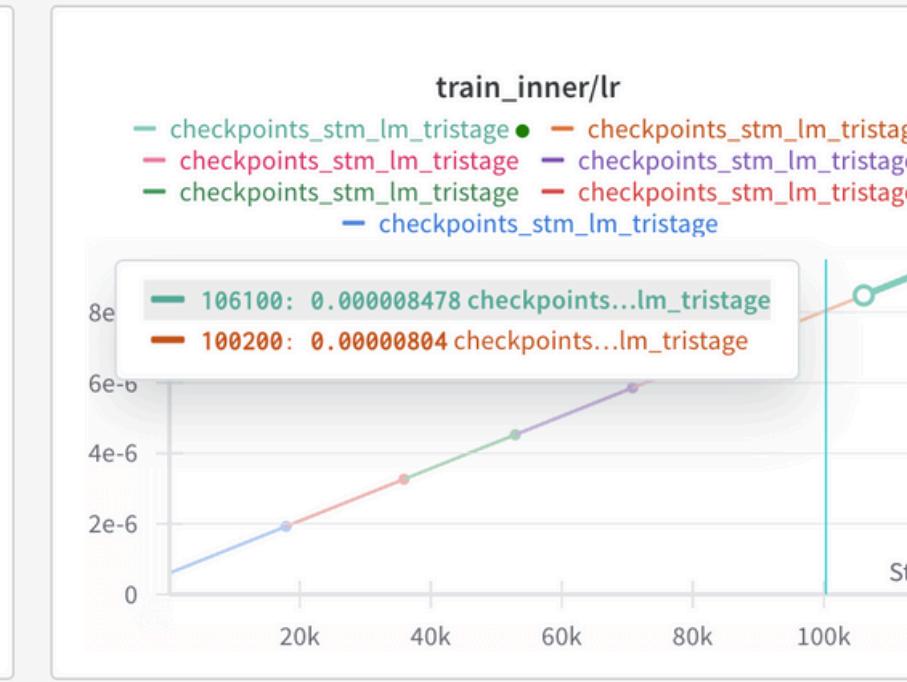
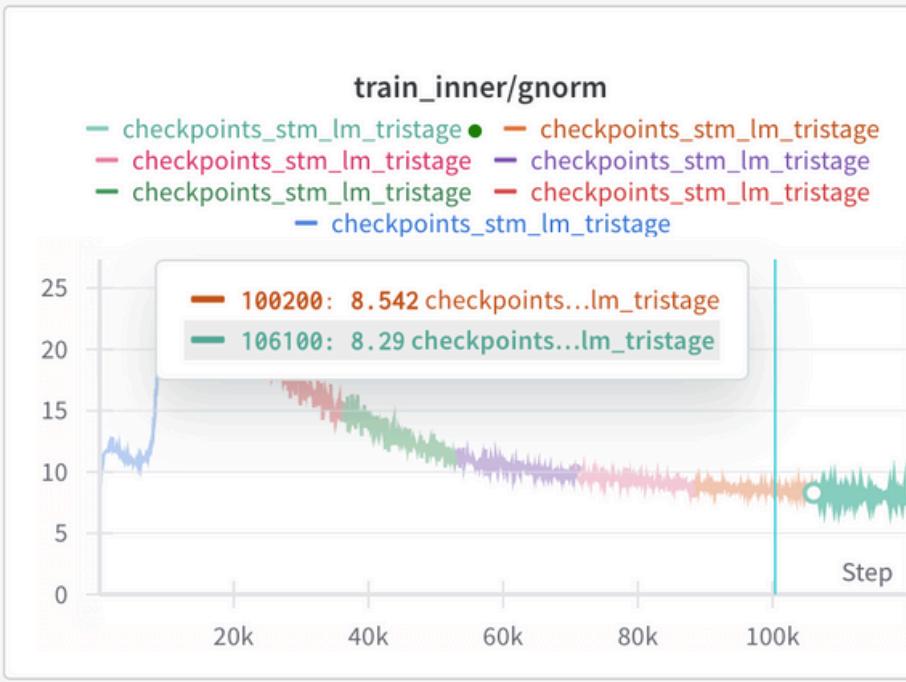
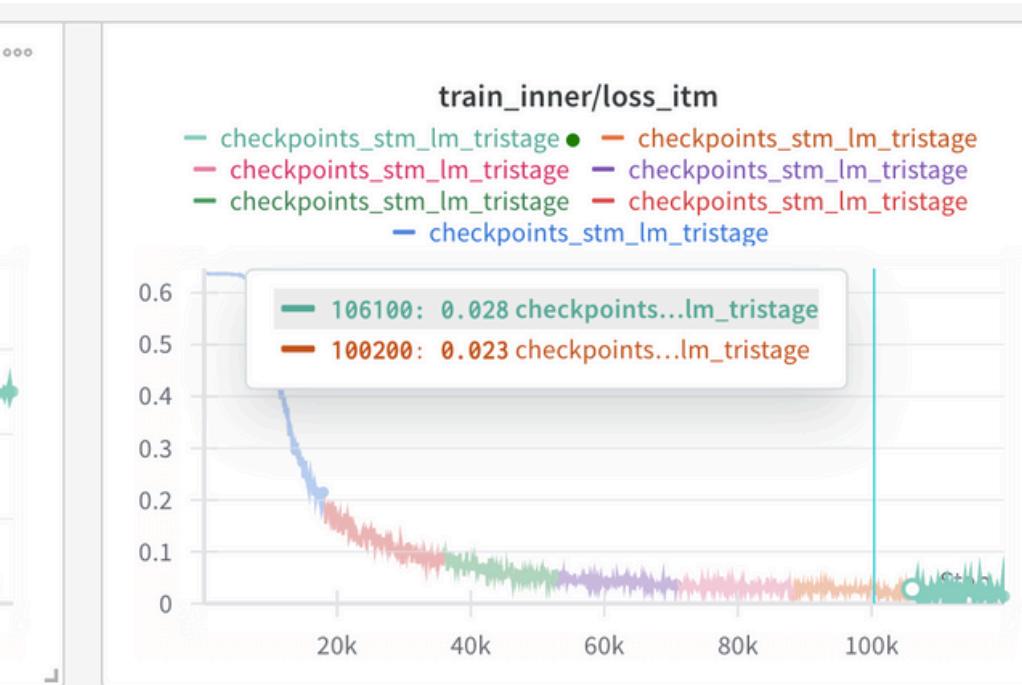
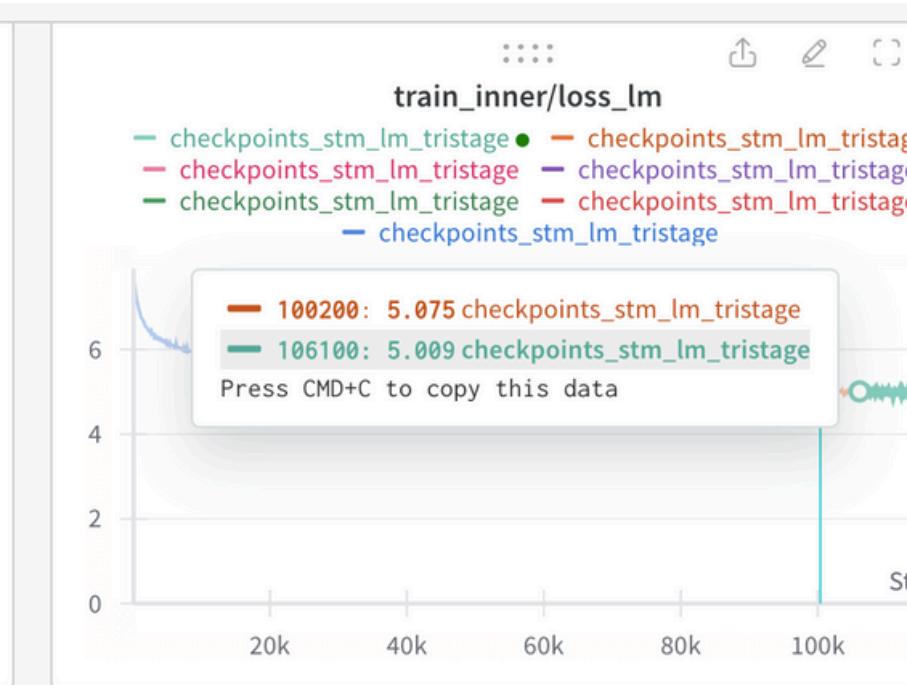
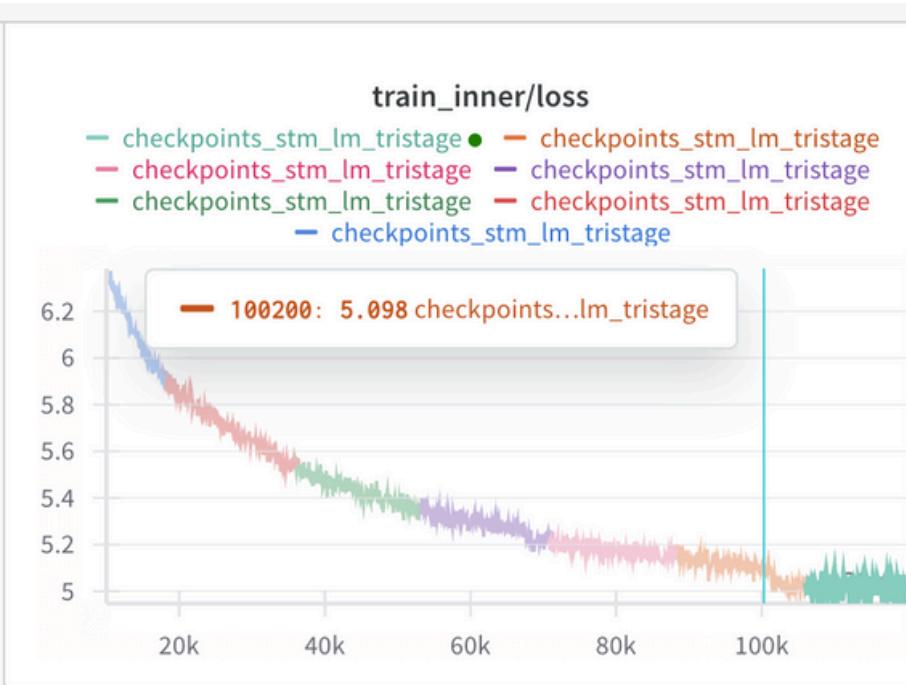
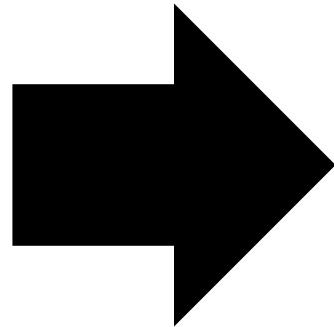
TRAINING OBJECTIVE

“Joint Speech-Text Matching & Speech-Text Generation Training”



$$\sum_{i=0}^N \underbrace{(-(y_i \cdot \log(p_i) + (1 - y_i) \cdot \log(1 - p)))}_{\text{Speech-Text Matching loss}}$$

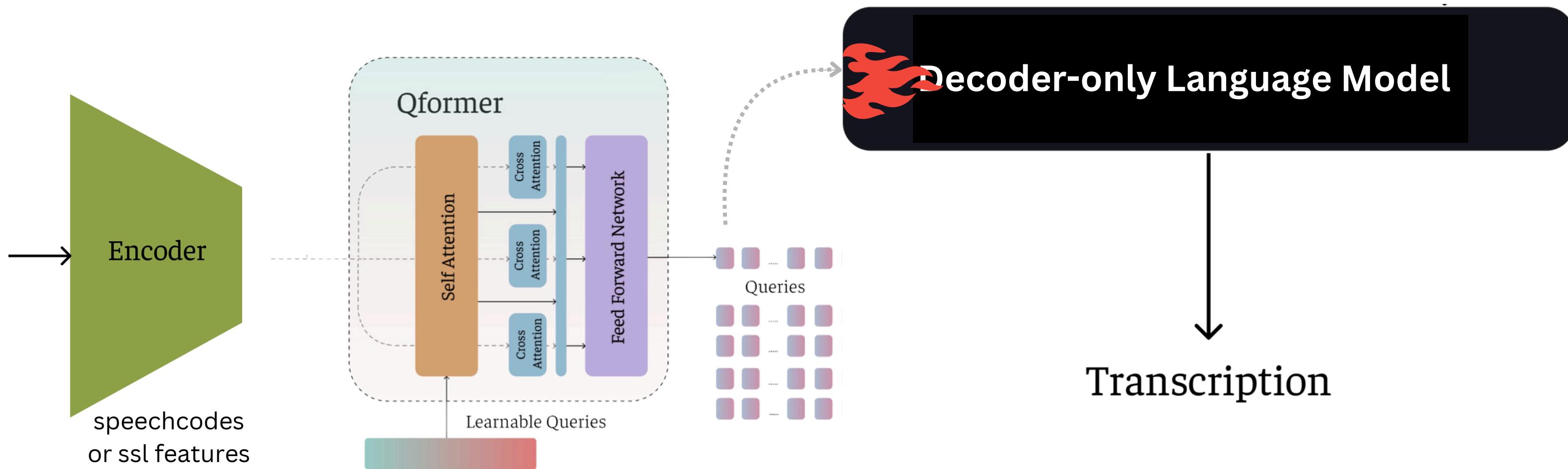
$$+ \sum_{i=1}^N \log (\mathbb{P}(u_i | u_{i-k}, \dots, u_{i-1}, \Theta))$$



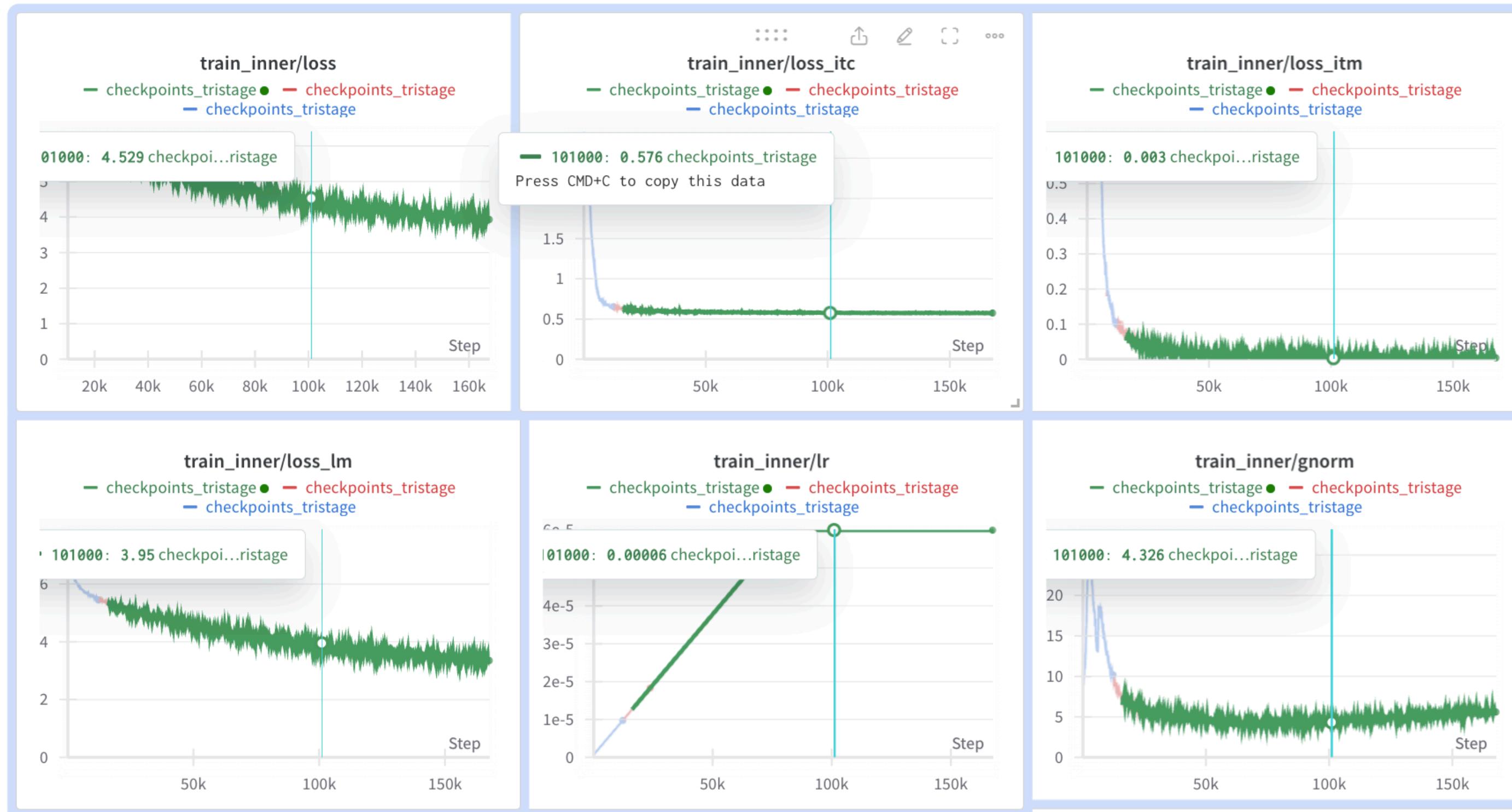
PROPOSED METHOD

detailed....

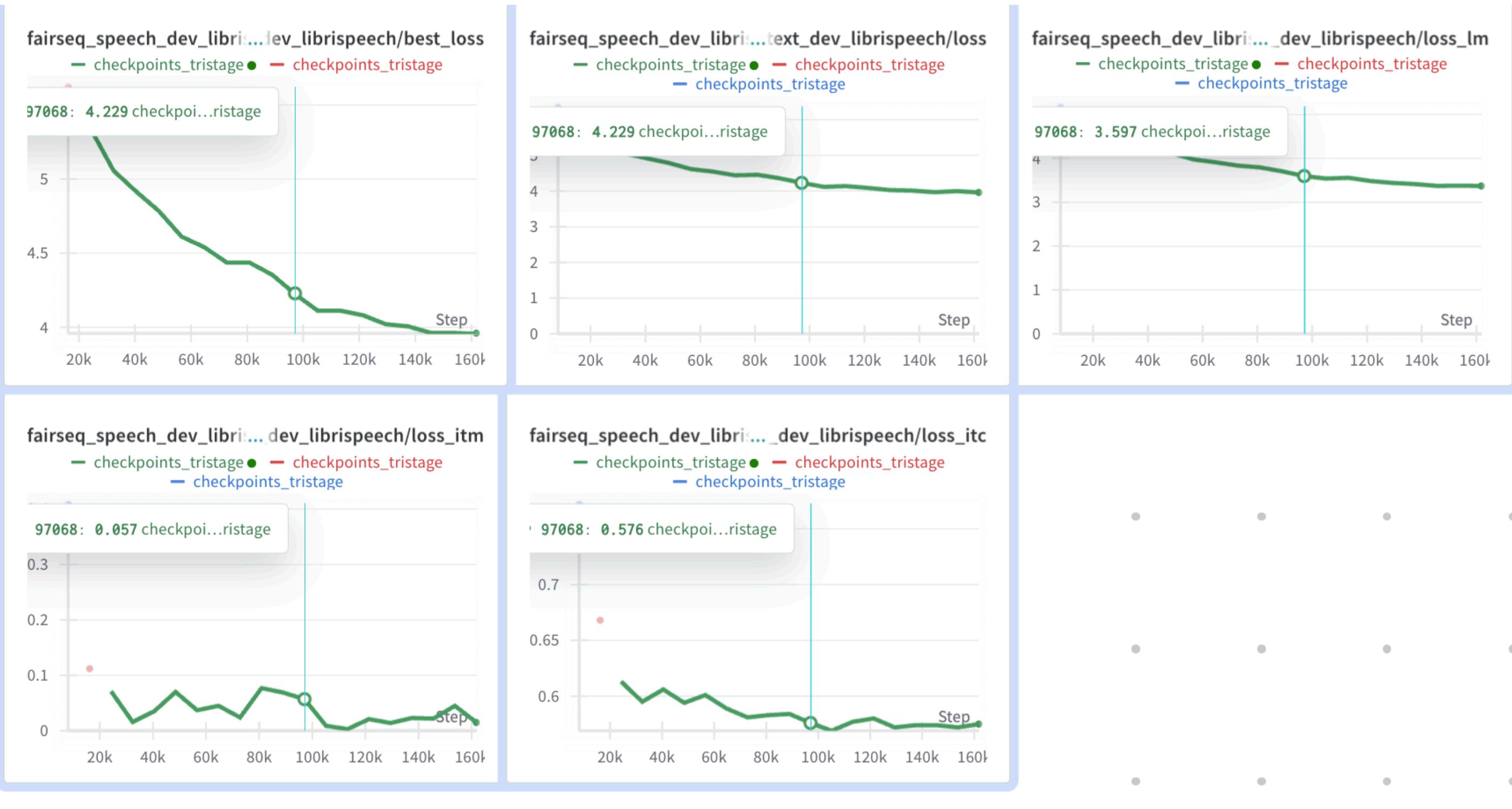
overall...



PROPOSED METHOD (TRAIN-STATS)

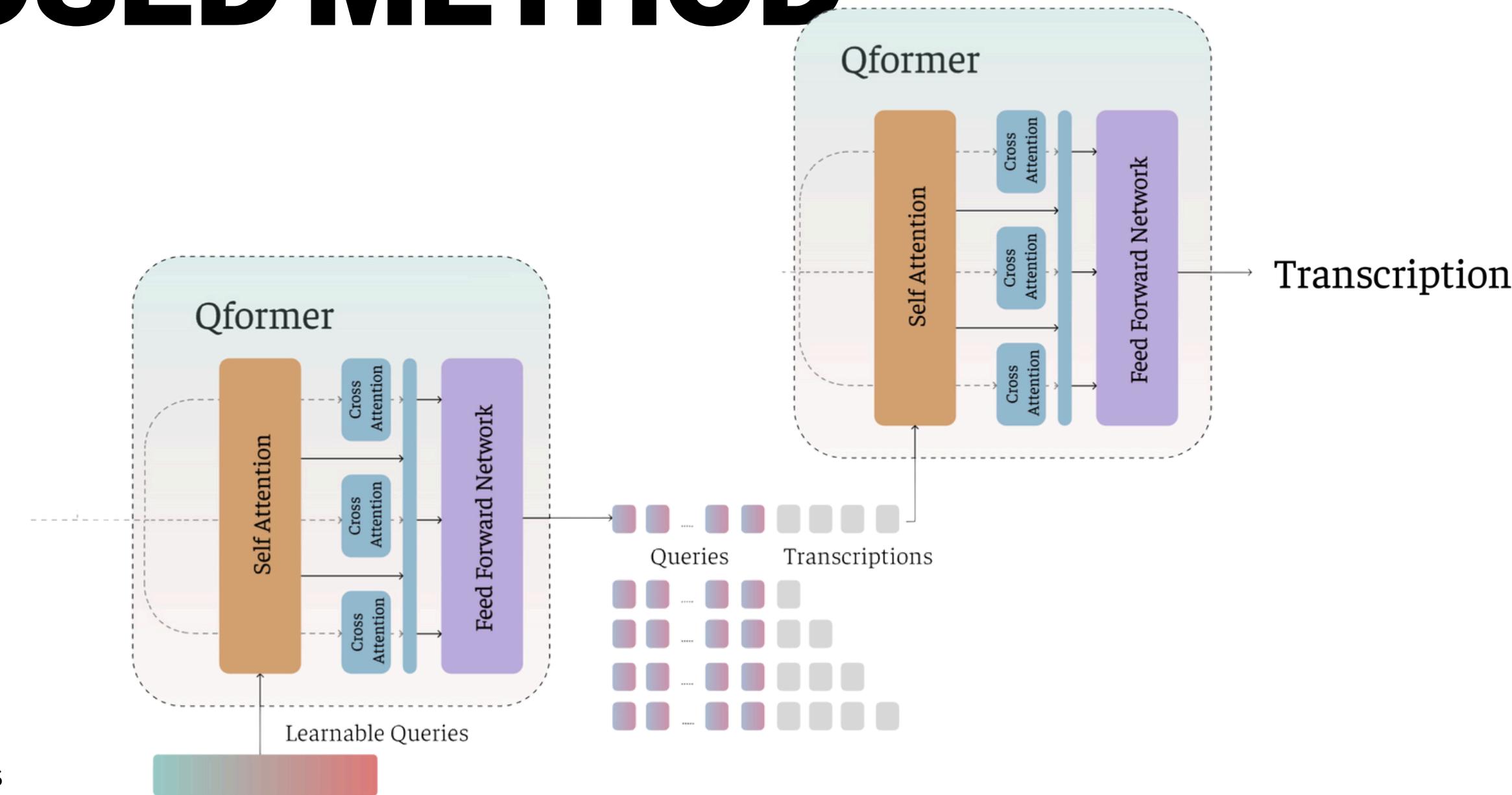
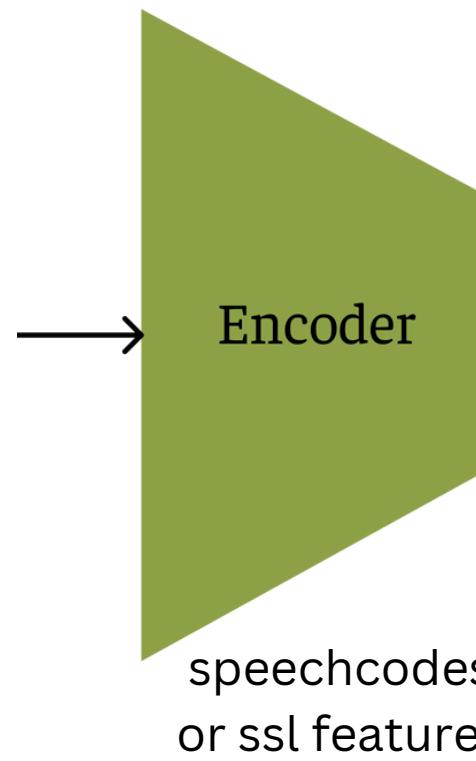


PROPOSED METHOD (DEV-STATS)



PROPOSED METHOD

real scenario



LIMITATIONS

The proposed system's connectivity was assessed based on its capability to generate corresponding text pairs not taking into account NLU specific datasets.

- **TLDR**, no NLU assessment

While an examination of attention heads and layers offered insights into the system's internal representations, we did not conduct any subsequent studies to delve deeper into these aspects.

- **TLDR**, no XAI, interpretations of attention

No ablation study examining training strategy



LLaMa 2



Vicuna



Mixtral

RESULTS

Experimental setup ASR-Qformer-*	Queries	WER↓		CER↓	
		dev-clean	dev-other	dev-clean	dev-other
SpeechTokenizer	50	103.56%	107.79%	79.15%	82.01%
	100	103.55%	109.16%	79.50%	83.54%
HuBERT-base	50	77.56%	77.56%	62.32%	62.32%
	100	74.80%	79.96%	60.70%	64.66%
HuBERT-large	50	52.61%	53.86%	44.02%	45.16%
	100	34.47%	39.61%	31.15%	32.84%

WER results are poor!!
BUT

RESULTS

Experimental setup ASR-Qformer-*	Queries	WER↓		CER↓	
		dev-clean	dev-other	dev-clean	dev-other
SpeechTokenizer	50	103.56%	107.79%	79.15%	82.01%
	100	103.55%	109.16%	79.50%	83.54%
HuBERT-base	50	77.56%	77.56%	62.32%	62.32%
	100	74.80%	79.96%	60.70%	64.66%
HuBERT-large	50	52.61%	53.86%	44.02%	45.16%
	100	34.47%	39.61%	31.15%	32.84%

WER results are poor!!
BUT

Experimental setup ASR-Qformer-*	Queries	Precision↑		Recall↑		F1↑	
		dev-clean	dev-other	dev-clean	dev-other	dev-clean	dev-other
SpeechTokenizer	50	84.74%	81.18%	84.83%	84.41%	84.78%	84.29%
	100	84.79%	84.06%	85.62%	84.58%	84.92%	84.31%
HuBERT-base	50	88.18%	87.10%	88.59%	87.90%	88.37%	87.49%
	100	88.46%	87.33%	88.87%	88.15%	88.65%	87.73%
HuBERT-large	50	92.11%	91.49%	92.32%	91.99%	91.73%	92.06%
	100	94.49%	93.56%	94.60%	93.83%	94.60%	93.83%

Semantic Similarity is better!!!

RESULTS

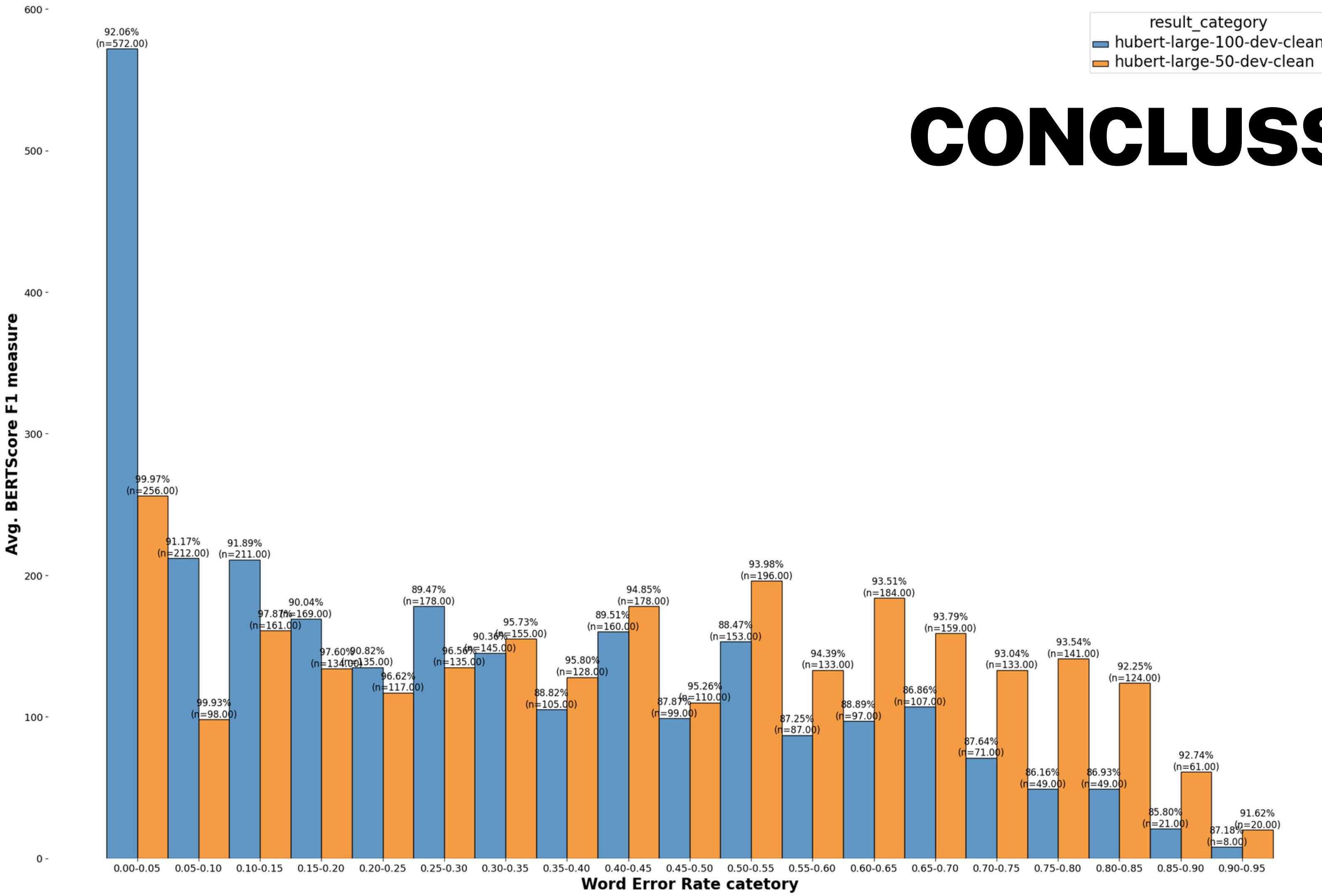
Experimental setup ASR-Qformer-*	Target Transcription	Predicted Transcription	WER
HuBERT-base-50	there is a seat in the garden at the side of the house again she hesitated	there is a seat in the garden at the side of the house again she hesitated	0.0%
	never had any act seemed so impossible	never had any accident seemed so impossible	14.28%
	no she shut her eyes	no she shut her eyes and i	40.00%
	then she turned towards the quarter indicated and disappeared round the laurel bushes	then she turned toward the quarter below town and ordered lord chelford the coroner's court	61.53%
	look here he said this is all nonsense you know you are tired out and there's something wrong what is it	look here he said this is all you know and there is something wrong here you are tired and tired	57.14%
HuBERT-large-50	i shall lock up all the doors and windows in the house and then i shall give you my latch key and you can let yourself in and stay the night here there is no one in the house	i shall lock up all the doors in the house and then i shall lock up your door and let you in there is no one in the house and you can stay all the night here myself and the lazaretto	61.53%
	i will catch the night train and bring my mother up to morrow then we will see what can be done	i will catch up the night train and bring my mother to morrow then we will see what can be done	9.52%
	when she opened the door of it the bright fire which beenie undesired had kindled there startled her the room looked unnatural uncanny because it was cheerful	when she opened the door of it the bright fire sparkled there unlike the room which linley had described the room was unmoved because it looked untenanted	51.85%

Experimental setup ASR-Qformer-*	Queries	WER↓		CER↓	
		dev-clean	dev-other	dev-clean	dev-other
SpeechTokenizer	50	103.56%	107.79%	79.15%	82.01%
	100	103.55%	109.16%	79.50%	83.54%
HuBERT-base	50	77.56%	77.56%	62.32%	62.32%
	100	74.80%	79.96%	60.70%	64.66%
HuBERT-large	50	52.61%	53.86%	44.02%	45.16%
	100	34.47%	39.61%	31.15%	32.84%

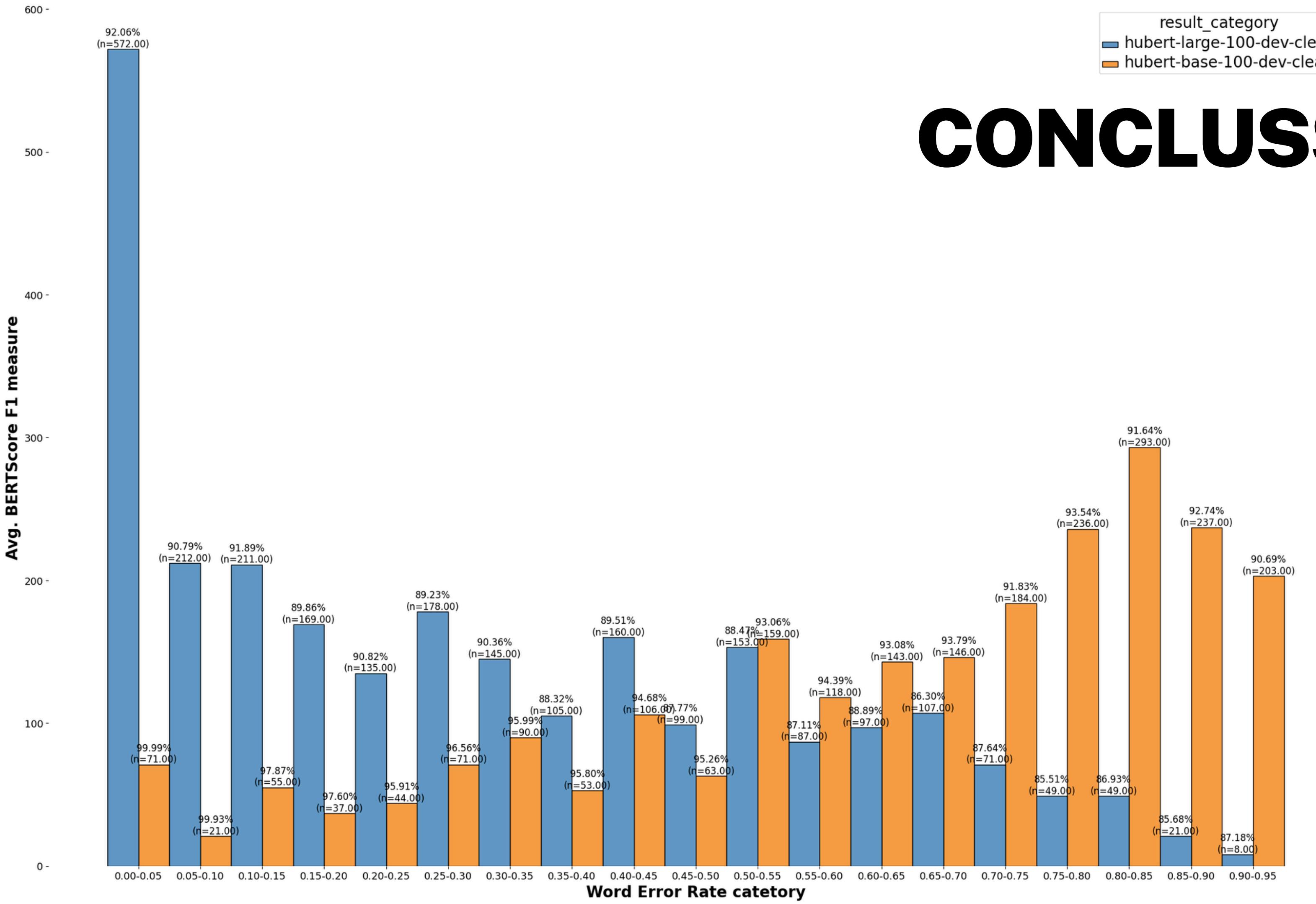
Experimental setup	WER↓		CER↓	
	dev-clean	dev-other	dev-clean	dev-other
CTC/AED FBank [48]	2.5%	6.3%	—	—
CTC/AED discrete tokens [48]	2.2%	4.5%	—	—
ASR-Qformer-HuBERT-large-100	34.47%	39.61%	31.15%	32.84%

Table 5.2: WER Comparison on LibriSpeech from existing methods

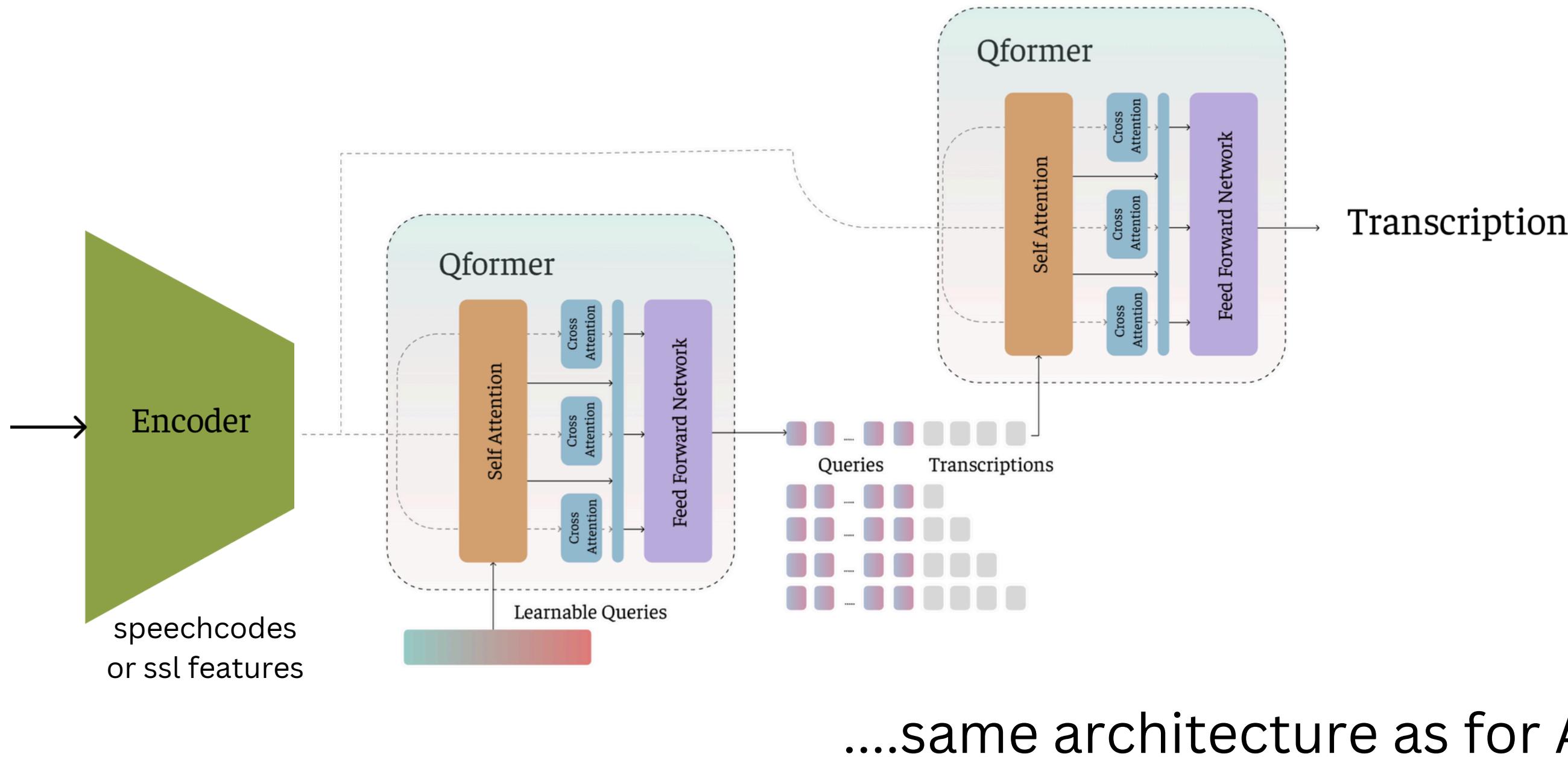
CONCLUSION



CONCLUSION



PROPOSED METHOD FOR AST



* AST - Automatic Speech Translation

RESULTS

System	BERTScore	
	En→De	En→Fr
STRONGBASELINE* [51]	77.44%	81.75%
WEAKBASELINE* [51]	74.86%	77.28%
APPTEK-Constrained [51]	77.32%	–
HW-TSC-Constrained [51]	–	76.11%
HW-TSC-Unconstrained [51]	75.79%	–
APV-Unconstrained [51]	73.68%	77.77%
HW-TSC-Constrained [51]	74.07%	76.11%
AST-Qformer-HuBERT-large-50	71.24%	65.39

Table 5.6: Combined En→De and En→Fr translation evaluation for BERTScore

Experimental setup	Precision↑		Recall↑		F1↑	
	dev	tst-COMMON	dev	tst-COMMON	dev	tst-COMMON
HuBERT-large En→De	71.35%	71.48%	71.20%	71.24%	71.25%	71.33%
HuBERT-large En→Fr	69.37%	70.54%	69.47%	70.27%	69.39%	70.37%

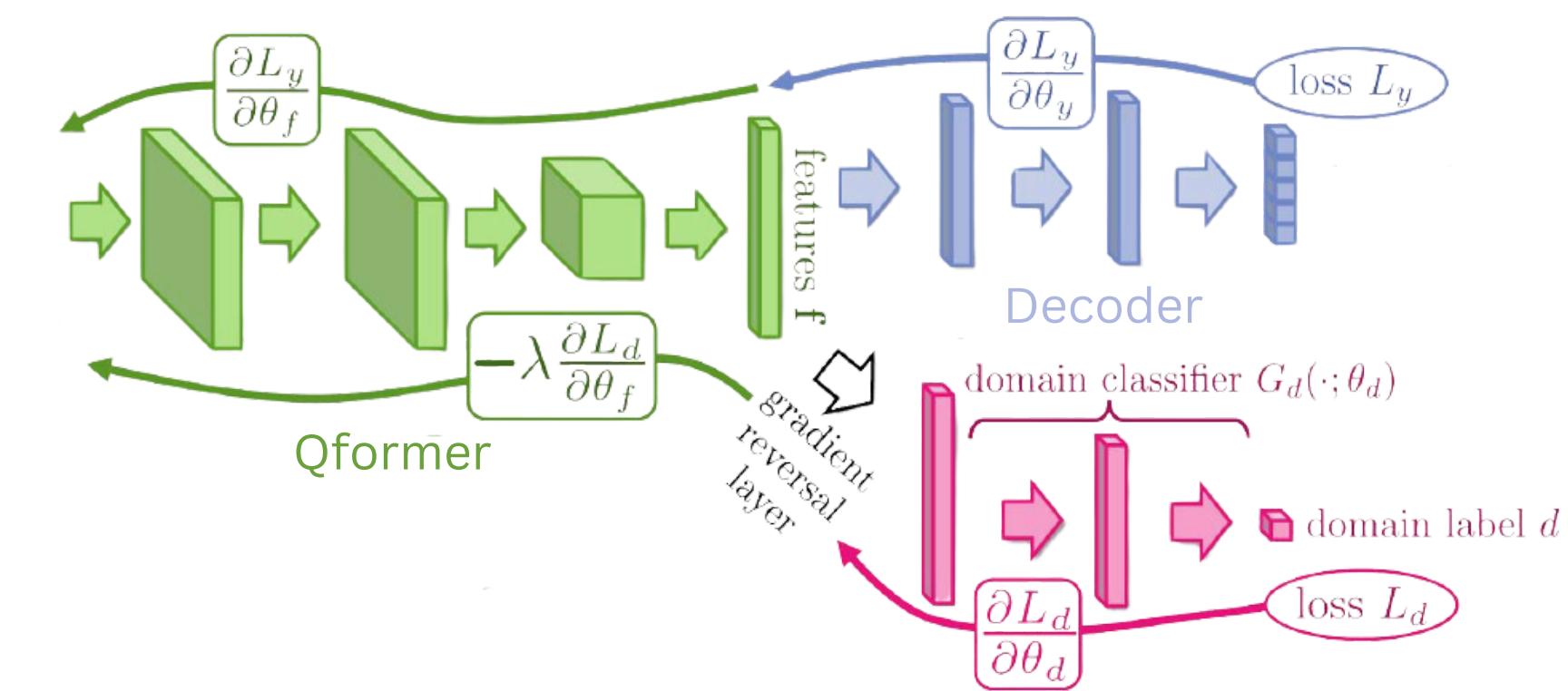
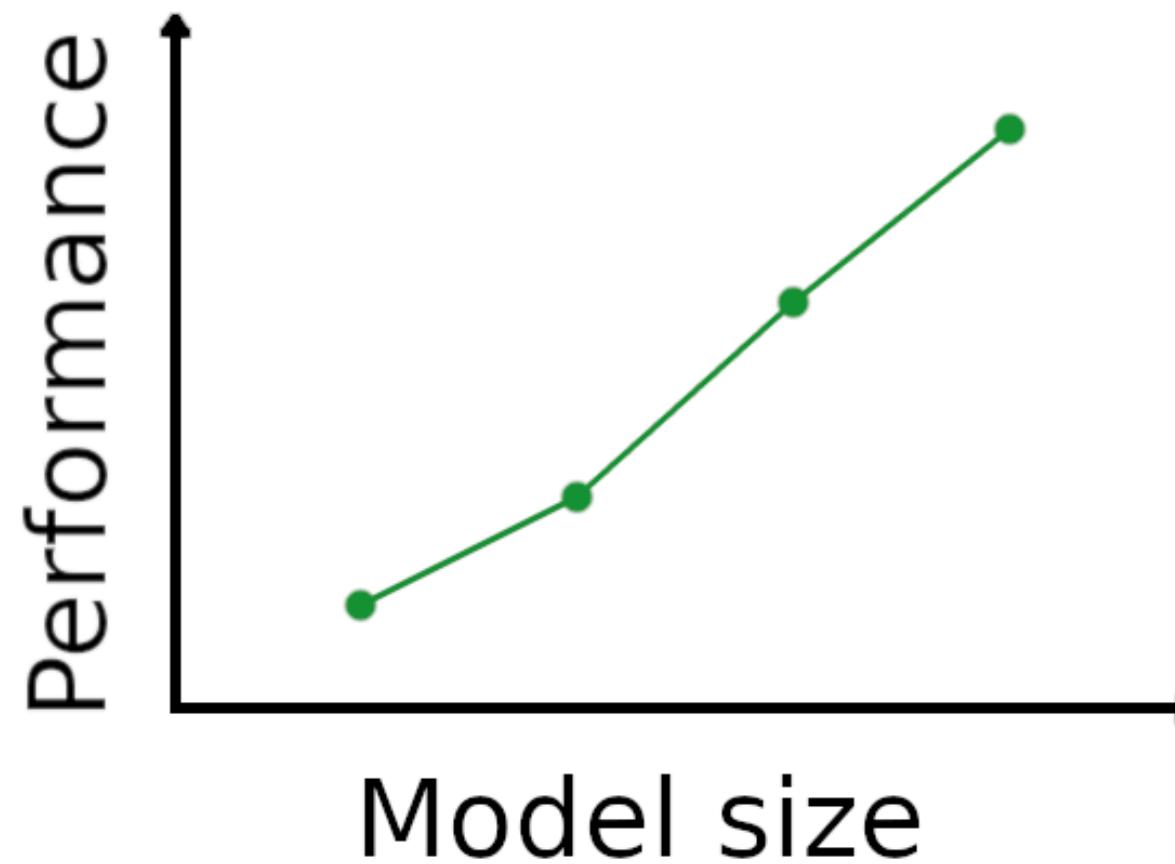
Table 5.7: BERTScore Comparison on LibriSpeech development and test splits

BLEU exhibits same problem as
WER & CER discussed
previously

... BERTScore

DISCUSSION

Insights from the drawn figures and tables emphasize the importance of query size and model depth as a critical factor in the design and optimization of speech recognition systems, highlighting its impact on the system's ability to process and interpret speech data effectively. Results from test-clean, test-other split from LibriSpeech indicate that the proposed method suffers distributional shifts raising the need for a domain adaptation strategy method. We believe that there is plenty of potential iterative work in the direction of ASR-Qformer in the provided research.



DISCUSSION

Normally, current ASR models operate with vocabulary on character level with classification objective, however, Qformer uses BERT's weights and vocabulary. Thus in our experiments we align speech representation with BERT's vocabulary that operates on word pieces rather than usual character level. This obviously adds complexity to the system.

Character level tokenization



Sub-word level tokenization



SpeechTokenizer uses CNN feature encoder/extractor and RVQ with 8 vector quantization block each with 1024 sized codebook

HuBERT used CNN feature encoder/extractor

example: 21seconds audio with sampling rate 16K is 336K frames

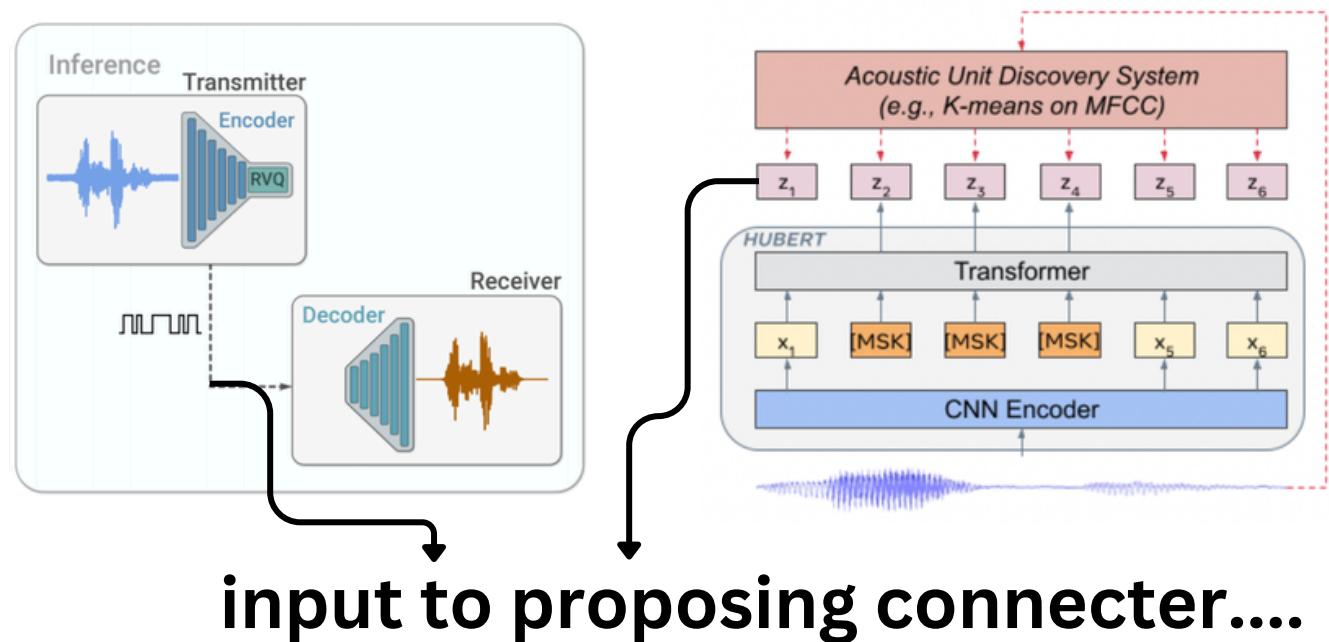
- SpeechTokenizer input tensor is **[1, 1050, 1024]**
- HuBERT input tensor is **[1, 1049, 768]**
- Tokenized text tensor is **[1, 75, 768]**

----->due to semantic knowledge distillation

DISCUSSION

The outcomes derived from the translation tasks further emphasize the potential for advancement within our architectural framework. The comparative analysis with benchmarks from the IWSLT 2022 Isometric Spoken Language Translation challenge illuminates the disparities in performance, particularly in semantic understanding and translation accuracy as reflected by our lower BERTScore. This revelation not only highlights the commendable capabilities of competing networks but also serves as a motivation for the refinement of our architecture.

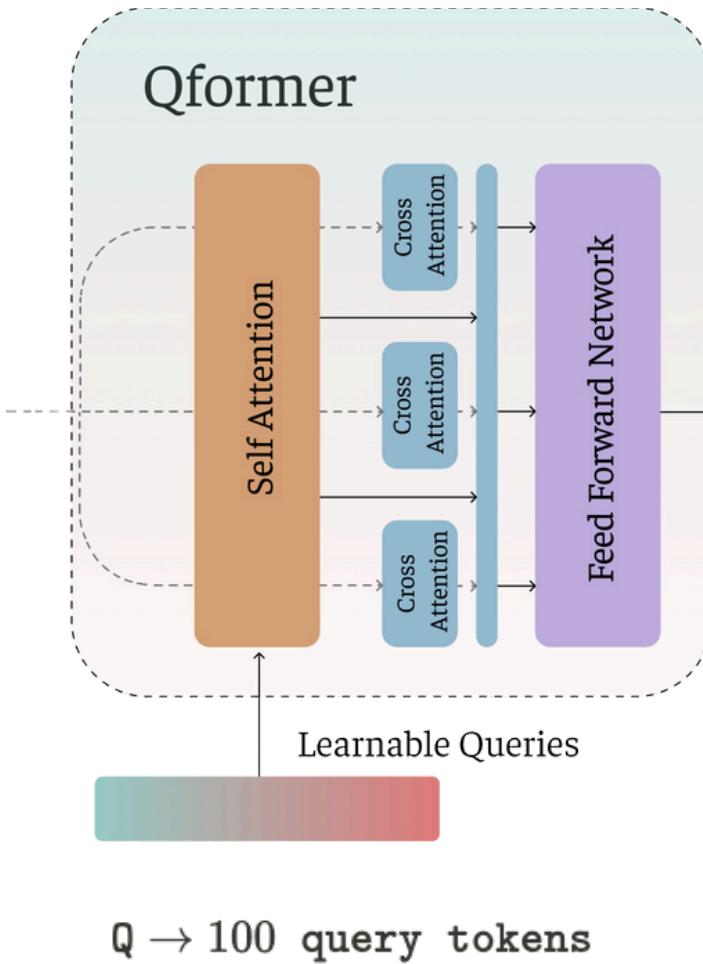
Overall, there is still a gap between speech and text modality, despite the main outcomes suggesting that the proposed training method is feasible and experiments representing semantic similarity to the target transcription. By that, we mean systems that connect LLM with speech encoders that were not created with speech-to-text objective. Moreover, positive insights are not revealed by the neural encoder's speech code results. The latent representation we transferred from a quantized codebook may be the cause of the subpar performance. Proposed hypothesis posited that the encoder's representations with hidden dimension were optimally suited as inputs for the Qformer, might be wrong, yielding the need in introducing embedding layer.



Suggestions...

1. Clustering HuBERT's intermediate layer representation
2. Introducing embedding layer for speechcodes from SpeechTokenizer rather than using codebook's latent representation

CONCLUSION



TLDR; “We introduce an alignment strategy inspired by the vision-language model BLIP2, adapting it by replacing speech representations. This approach does not iterate upon BLIP's original methodology but rather expands the application of discrete speech representation. Our study investigates the effectiveness of using speech units as inputs across various speech processing tasks, such as Automatic Speech Recognition and Automatic Speech Translation. To assess the versatility of discrete units, we utilized two distinct types of speech encoders: one focused on reconstruction and the other on a masking objective.

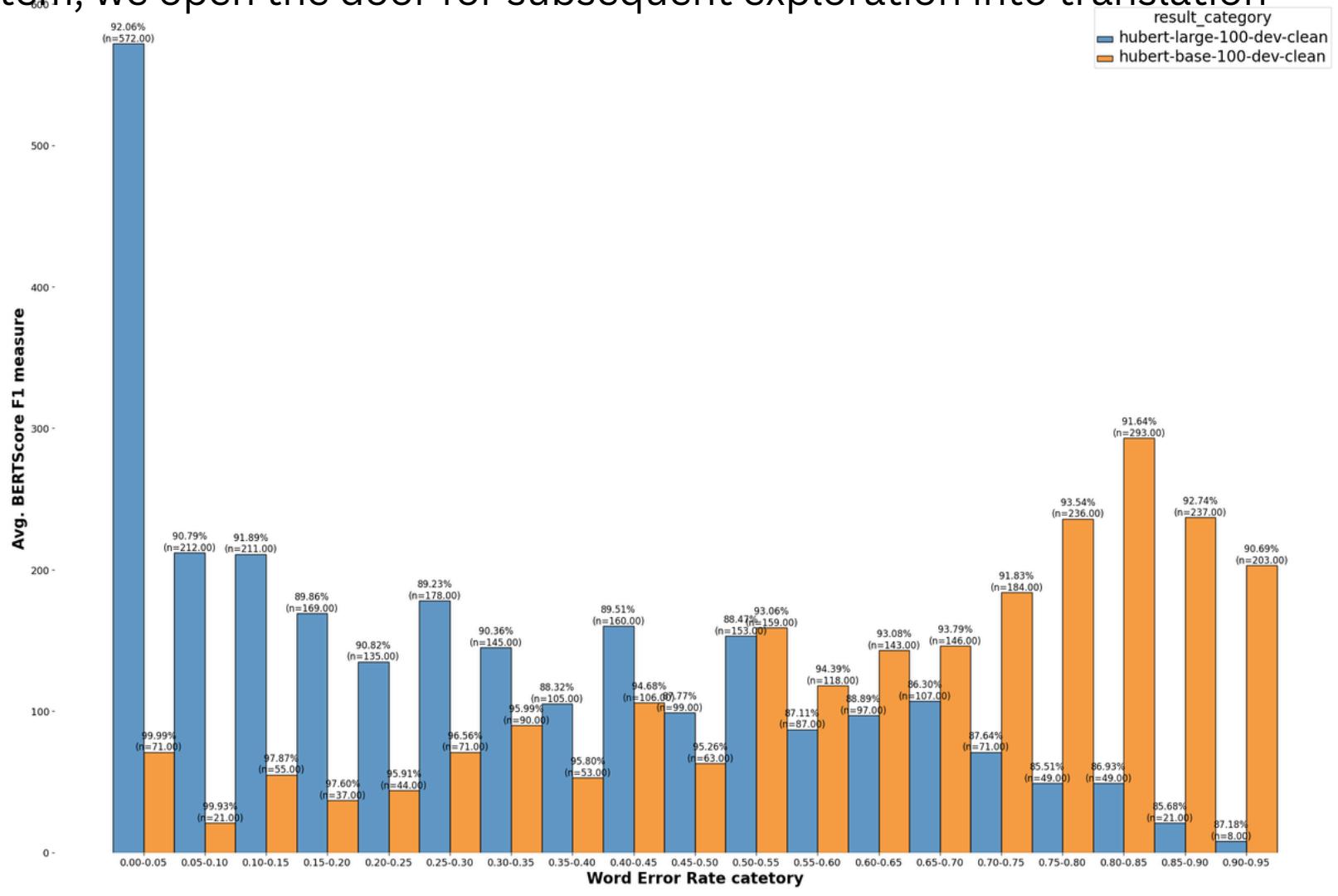
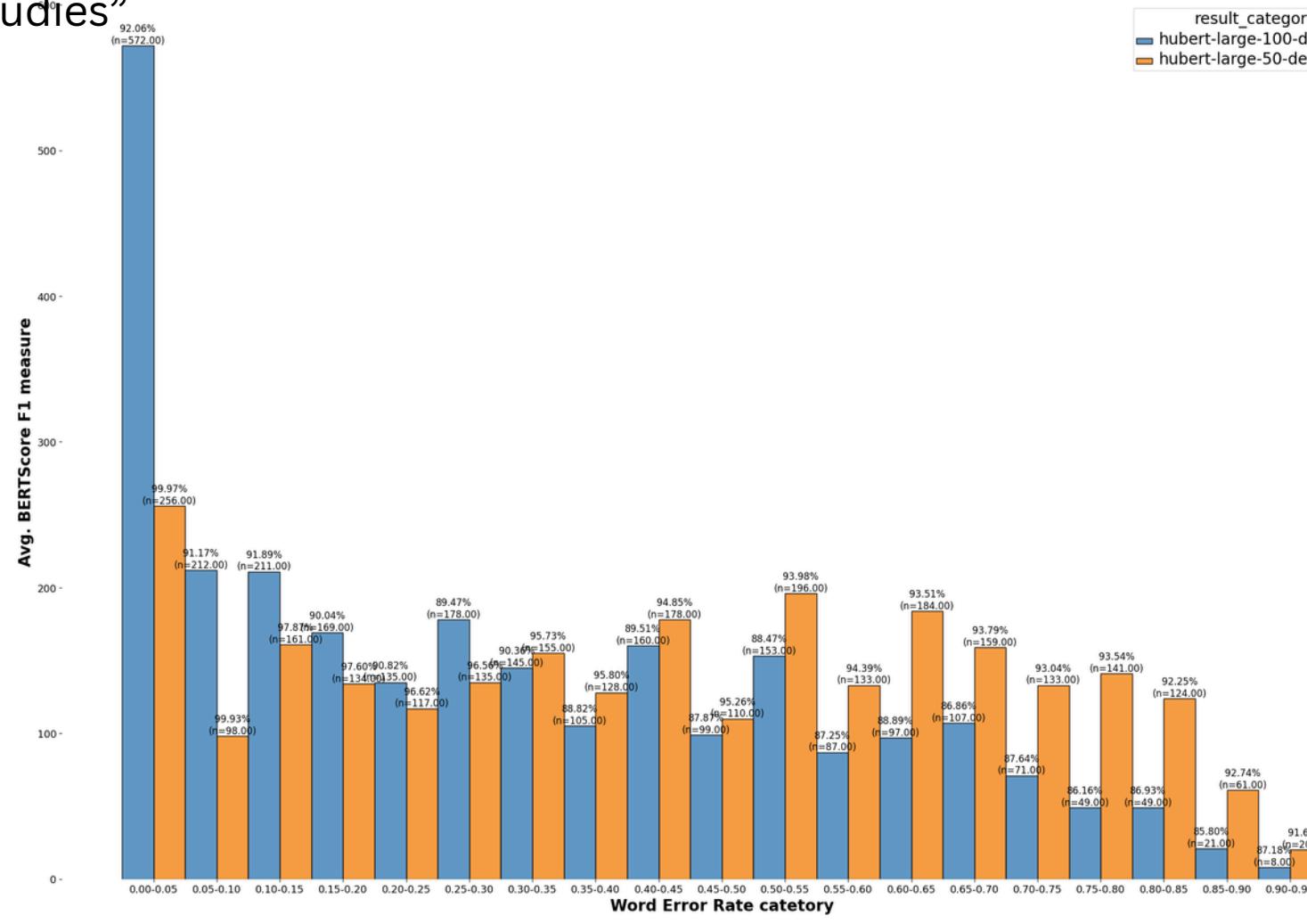
Our experiments conducted on the LibriSpeech dataset demonstrate the efficacy of our method, though it does present challenges in precisely matching the target text. By comparing our model's performance to existing benchmarks, we found a interrelationship between the outputs and the semantic representation of the gold standard labels. However, the neural encoder designed for speech reconstruction, which we employed, did not produce favorable results based on our assumptions, thereby suggesting that features derived from the masking model offer a more suitable context for extracting queries. The research questions that we outlined in the beginning of the paper raise two main problems that were addressed in this study.”

NLPR; (not long pleaser read)

- Introduced alignment strategy tested on ASR and AST
- Speechcodes and SSL representation is tested as input
- LibriSpeech results suggest that exact word matching did not succeed but semantic similarity is preserved

CONCLUSION

TLDR; “Firstly, our investigation into whether speech-encoded representations, derived either from compression algorithms or speech self-supervised learning methods, can bridge the modality gap between text and speech concludes positively. Evidence from our ASR-HuBERT experiments supports the efficacy of larger encoders with an increased number of query tokens. A more detailed analysis provided in our evaluation section’s Figure\ref{fig:hubert_base_vs_large_100} depicts that, across various error categories represented on the x-axis, semantic meaning is preserved with approximately 90% similarity. This underscores the ability of self-supervised learning features to align closely with text. The results of our findings are particularly promising for future research focused on low-resource languages with substantial amounts of untranscribed speech resources. By demonstrating that training a speech model without actual transcriptions and training a text model with a modest quantity of speech-text pairs is sufficient to construct an ASR system, we open the door for subsequent exploration into translation studies”



CONCLUSION

TLDR; “Finally, this study contributes a novel reduction strategy through the use of a fixed number of query tokens to extract information from the speech modality. This approach addresses a common challenge within the speech research community: the search for compressing strategy in high-dimensional speech representations into more manageable, lower-dimensional forms to reduce computational resource and storage requirements. Our methodology effectively condenses speech representations into a fixed number of tokens that approximate the dimensionality of text tokenized representations, offering a practical solution to these challenges.”

NLPR

- Another advantage of the system is fixed-sized reduced speech representation that might be used for autoregressive decoding

