

Towards Post-training Alignment for Large Speech Models

Speaker: Hu Yuchen (NTU 4th-year Ph.D. Student)

Date: Sep 30, 2024



Background: Large Speech Models

Whisper

SeamlessM4T

Canary

VALL-E

VoiceCraft

CosyVoice

SpeechGPT

Qwen-Audio

SALMONN

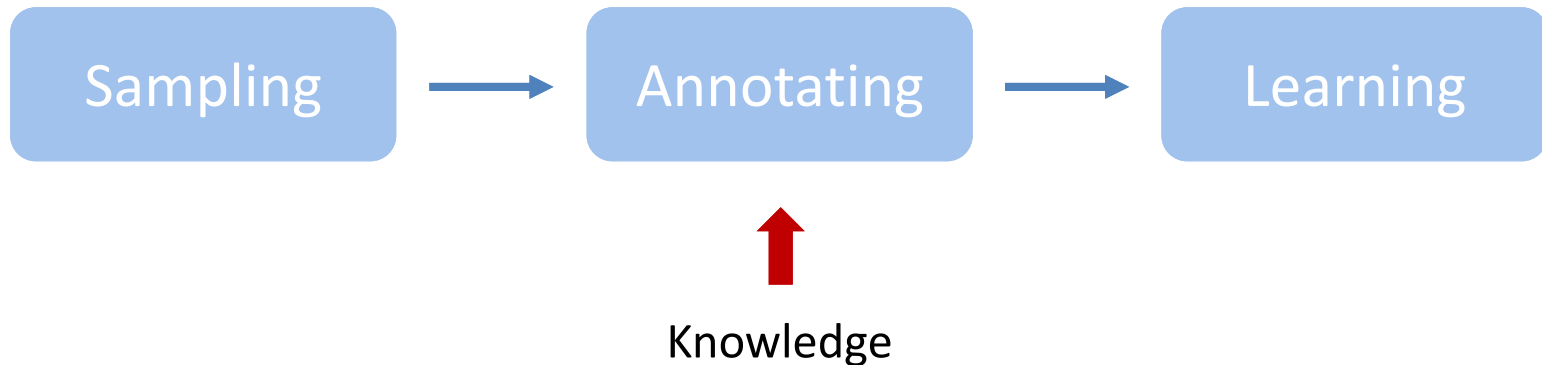
How to **efficiently** align them to specific scenarios?

Background

Pretraining & SFT:

$$\mathcal{L}(x, y) = \sum_{l=1}^L -\log \mathcal{P}_{\theta}(y_l | y_{l-1}, \dots, y_1, x) \quad \text{Learn a distribution --> uncertain}$$

Post-training:



Outline

Speech-to-Text Understanding

- Motivation
- Method
- Results on ASR & AST
- Conclusion & Discussion

Text-to-Speech Synthesis

- Motivation
- Method: UNO
- Experimental Results
- Extension: RIO
- Conclusion & Discussion

Self-Taught Recognizer: Toward Unsupervised Adaptation for Speech Foundation Models

Yuchen Hu^{1*}, Chen Chen^{1*}, Chao-Han Huck Yang², Chengwei Qin¹,
Pin-Yu Chen³, Eng Siong Chng¹, Chao Zhang⁴

¹ Nanyang Technological University ² NVIDIA Research

³ IBM Research ⁴ Tsinghua University



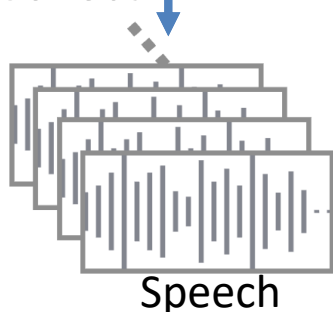
Motivation

To deploy an ASR system in a practical scenario:



(NTU Canteen)

Collect

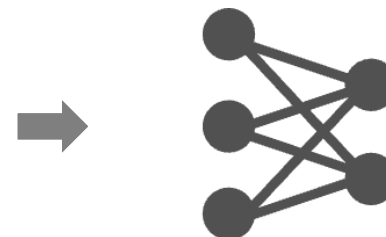


Manual
Labeling
+



Transcription

Supervised Domain Adaptation

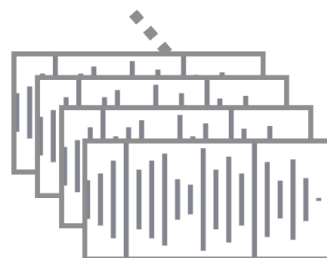


Neural Model

A very convenient approach is:



+



Unlabeled
Speech

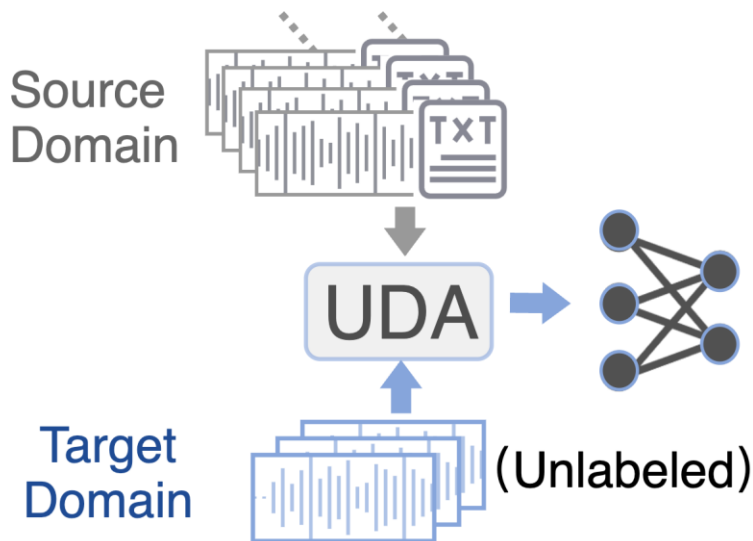
Unsupervised
Domain Adaptation



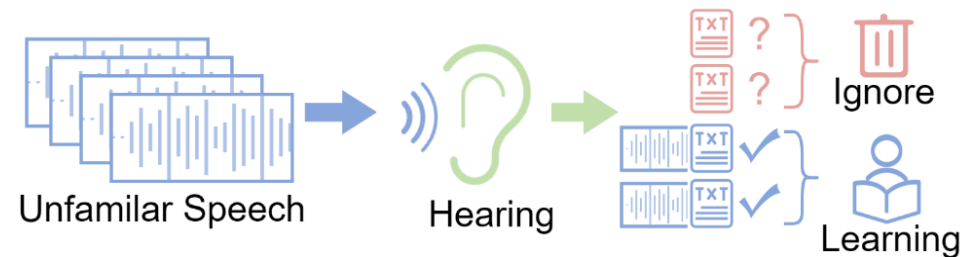
UDA

Motivation

UDA in ASR:



Human's UDA solution:

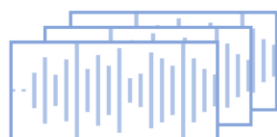


- “Unsupervised ” is for adaptation process, but the learning schedule is semi-supervised.
- Considering the exhibited ability of large speech model:
Can we **skip the source-domain data** for target domain adaptation? ➔ Source-free UDA

Method (Self-training --> Post-training)

1) Pseudo Labeling:

Source-Free



(Sampling)



Pseudo Label



Keep this utterance or discard?

--> Monte Carlo sampling

2) Informed Finetuning:



How to assign weight for **each token**?

Unlabeled Speech



(Learning)

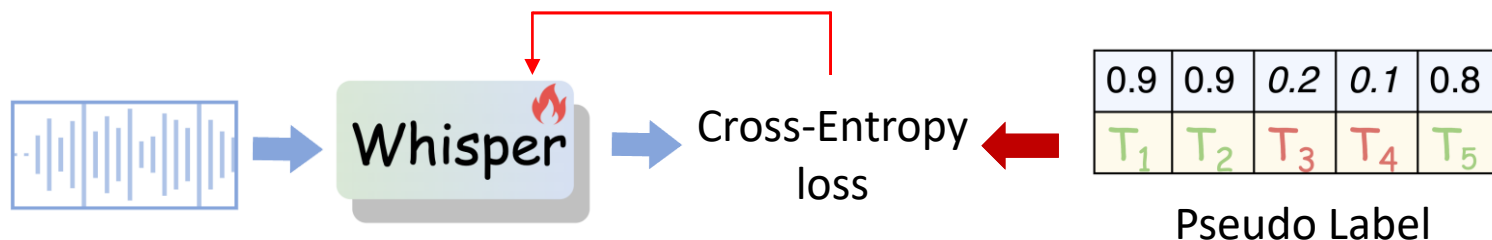
Cross-Entropy loss

(Annotating)

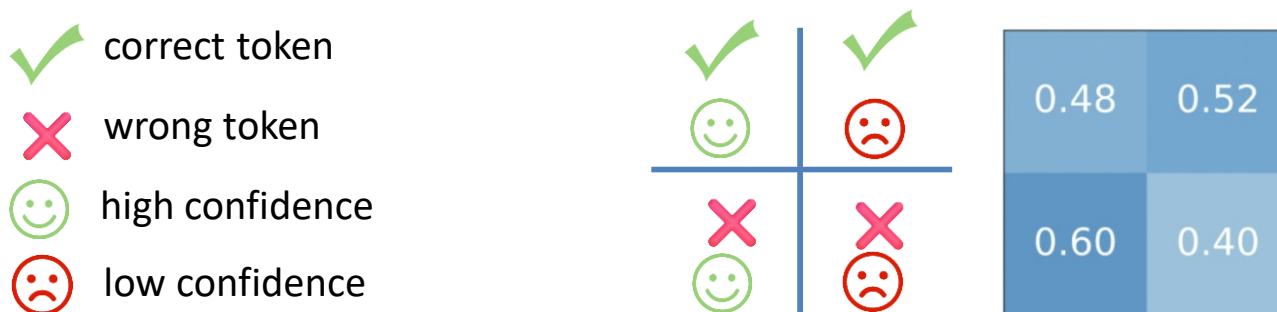
0.9	0.9	0.2	0.1	0.8
T ₁	T ₂	T ₃	T ₄	T ₅

Pseudo Label

Candidate 1: Confidence Score



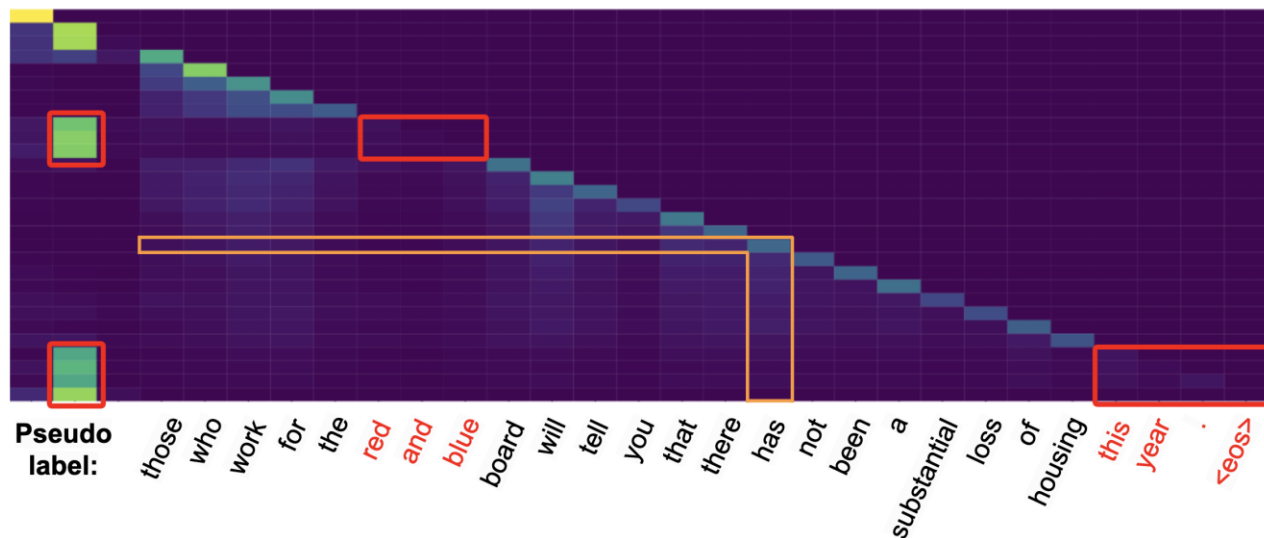
Experimental observation: decoding performance on CHiME-4 test-real



Confidence score is unreliable!

Candidate 2: Self-Attention Matrix

<|transcribe|>



Attentive score:

$$\mathcal{A}_l = \sum_{j=4}^l W_{l,j} + \sum_{i=l+1}^L W_{i,l},$$

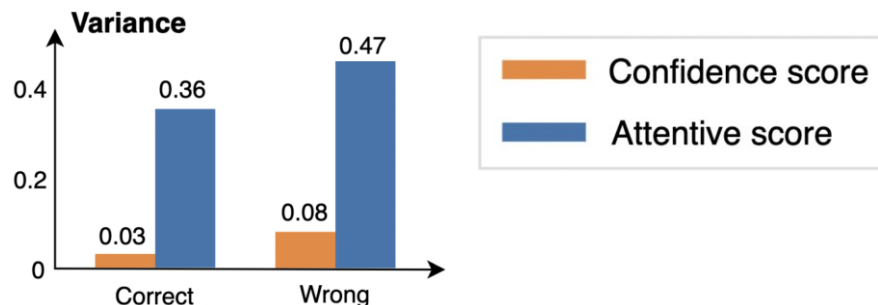
The importance of l -th token in whole utterance^[8]

Correct/**Wrong**

Is A_l more **reliable** than C_l ?



Is A_l stable for guide finetuning?



Conclusion: attentive score is **more reliable** but **less stable** than confidence score.

10



STAR: Integrate A and C for each token

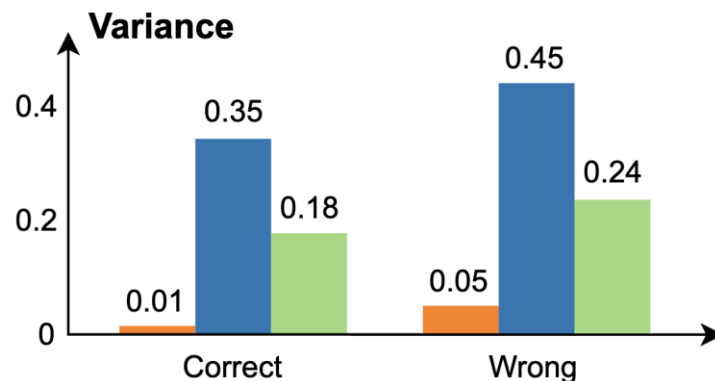
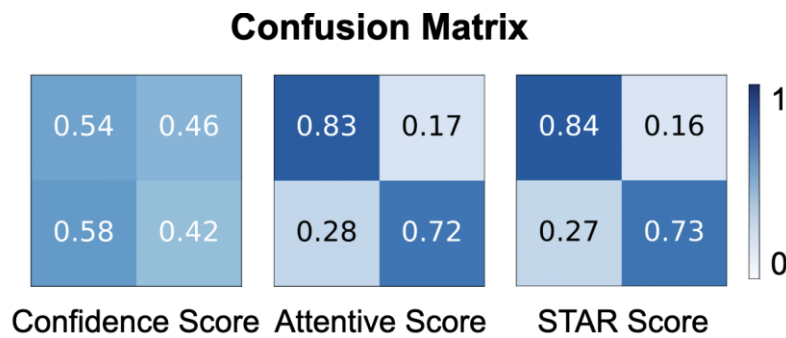
Criteria: - If A-C conflict, then follow A:

$$\mathcal{S}_l^{\text{conf}} = [\sigma(\mathcal{A}_l^2/\mathcal{C}_l - \lambda) + \sigma(\mathcal{C}_l^2/\mathcal{A}_l - \lambda)] * \mathcal{A}_l$$

- If A-C consistent, then calibrate A using C:

$$\mathcal{S}_l^{\text{cons}} = [\sigma(\lambda - \mathcal{A}_l^2/\mathcal{C}_l) * \sigma(\lambda - \mathcal{C}_l^2/\mathcal{A}_l)] * \mathcal{A}_l * e^{(\mathcal{C}_l - \mathcal{A}_l)/\tau}.$$

Quick validation:



Effectiveness on Various Domains

STAR = Self-TAught Recognizer

Testing Scenario		Whisper (frozen)	Whisper (self-train.)	UTT _{filter}	TOK _{\mathcal{C}_l}	reweight \mathcal{A}_l	STAR (ours)	Whisper (real label)
		<i>Background Noise</i>						
CHiME-4	<i>test-real</i>	6.8	6.9	6.4	6.5	6.2	6.0 —11.8%	5.2
	<i>test-simu</i>	9.9	10.1	9.7	9.8	9.5	9.4 —5.1%	8.7
	<i>dev-real</i>	4.6	4.5	4.3	4.3	4.1	3.9 —15.2%	3.2
	<i>dev-simu</i>	7.0	7.0	6.6	6.7	6.6	6.4 —8.6%	5.9
LS-FreeSound	<i>babble</i>	40.2	37.6	35.0	33.5	31.3	30.2 —24.9%	27.2
	<i>airport</i>	15.6	15.5	15.2	15.3	15.0	14.8 —5.1%	14.5
	<i>car</i>	2.9	3.0	2.8	2.8	2.6	2.5 —13.8%	2.4
RATS	<i>radio</i>	46.9	47.2	46.0	45.5	44.9	44.6 —4.9%	38.6
		<i>Speaker Accents</i>						
CommonVoice	<i>African</i>	6.0	5.8	5.5	5.4	5.0	4.8 —20.0%	4.6
	<i>Australian</i>	5.8	5.7	5.6	5.5	5.2	5.1 —12.1%	4.3
	<i>Indian</i>	6.6	6.5	6.3	6.4	6.1	6.0 —9.1%	5.7
	<i>Singaporean</i>	6.5	6.2	5.8	5.8	5.4	5.1 —21.5%	4.9
		<i>Specific Scenarios</i>						
TED-LIUM 3	<i>TED talks</i>	5.2	4.9	4.7	4.8	4.3	4.1 —21.2%	3.6
SwitchBoard	<i>telephone</i>	20.8	20.5	19.8	19.3	18.6	18.1 —13.0%	15.3
LRS2	<i>BBC talks</i>	8.5	8.3	7.6	7.9	7.4	7.0 —17.6%	5.6
ATIS	<i>airline info.</i>	3.6	3.5	3.3	3.3	3.2	2.9 —19.4%	2.0
CORAAL	<i>interview</i>	21.5	21.3	20.8	20.7	20.4	20.1 —6.5%	17.9

Whisper
zero-shot

Previous
Semi-ASR

Ours

Real-label
training

Analysis

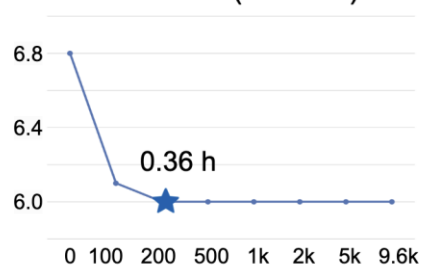
STAR can avoid forgetting:

Model	LS-FreeSound			RATS	CommonVoice				TED-3	SWBD	ATIS
	<i>babble</i>	<i>airport</i>	<i>car</i>		<i>af</i>	<i>au</i>	<i>in</i>	<i>sg</i>			
Frozen	40.2	15.6	2.9	46.9	6.0	5.8	6.6	6.5	5.2	13.3	3.6
Self-train.	38.2	16.6	2.9	47.3	6.4	5.9	6.7	6.3	5.3	13.7	3.4
STAR	33.3	15.7	2.8	46.1	6.1	5.8	6.7	5.6	5.0	13.5	2.9

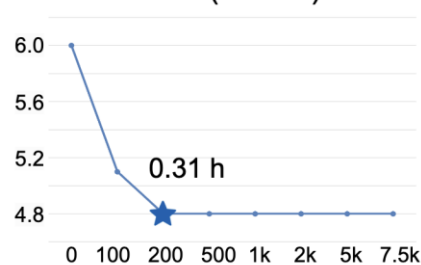
Train on CHiME-4 and test on OOD

STAR enjoys high data efficiency:

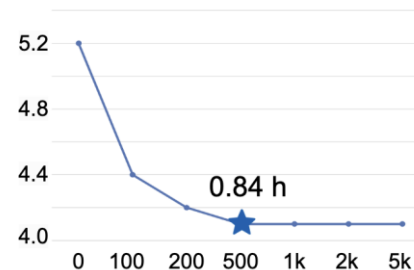
WER CHiME-4 (test-real)



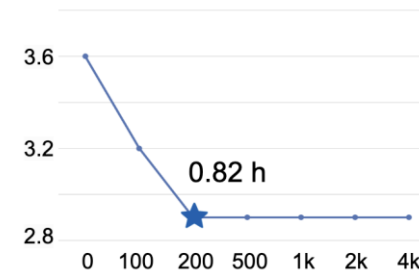
CV (African)



TED-LIUM 3



ATIS



train samples

Generalization

- Other models

Model	Baseline	Self-train.	STAR	Real
Whisper-V3-1.5B	6.8	6.9	6.0 _{-11.8%}	5.2
Whisper-Med-0.8B	8.9	8.8	8.0 _{-10.1%}	7.1
OWSM-V3.1-1.0B	8.4	8.1	7.5 _{-10.7%}	6.5
Canary-1.0B	8.2	8.0	7.2 _{-12.2%}	6.4
Parakeet-TDT-1.1B	8.0	7.8	7.0 _{-12.5%}	6.2

- Other task (Speech Translation on FLEURS)

X → En	Baseline	Self-train.	STAR	Real
Ar	21.9	22.1	23.3 _{+1.4}	24.5
De	33.7	34.0	35.9 _{+2.2}	36.5
Es	23.9	24.1	24.8 _{+0.9}	26.4
Fa	16.6	16.3	17.6 _{+1.0}	19.0
Hi	22.4	22.5	23.4 _{+1.0}	24.4
Zh	16.3	16.3	17.1 _{+0.8}	17.9

Ablation Study

- Different whisper sizes

Model Size	# Param.	Baseline	STAR	Real
large-v3	1,550 M	6.8	6.0 _{-11.8%}	5.2
large-v2		7.7	6.9 _{-10.4%}	6.0
large		7.5	7.0 _{-6.7%}	6.8
medium.en	769 M	8.9	8.0 _{-10.1%}	7.1
small.en	244 M	12.7	10.6 _{-16.5%}	9.0
base.en	74 M	32.4	17.7 _{-45.4%}	16.1

- Different training methods

Approach	# Param.*	Baseline	STAR	Real
Regular Finetuning				
Full	1550 M	6.8	6.0 _{−11.8%}	5.2
Enc-only	635 M		6.3 _{−7.4%}	5.0
Dec-only	907 M		6.1 _{−10.3%}	4.4
Parameter-Efficient Finetuning				
LoRA	16 M	6.8	6.0 _{−11.8%}	5.1
Reprogram.	0.4 M		6.7 _{−1.5%}	6.7

Iterative Post-training

Model	Test set	# Iterations						Real label
		0	1	2	3	4	5	
large-v3	<i>test-real</i>	6.8	6.0	5.9	5.7	5.7	5.7	5.2
medium.en		8.9	8.0	7.9	7.9	7.8	7.8	7.1
small.en		12.7	10.6	10.3	10.3	10.3	10.3	9.0
base.en		34.4	17.7	17.2	17.2	17.0	17.0	16.1
large-v3	<i>test-simu</i>	9.9	9.4	9.3	9.0	8.9	8.9	8.7
	<i>dev-real</i>	4.6	3.9	3.9	3.8	3.8	3.8	3.2
	<i>dev-simu</i>	7.0	6.4	6.4	6.4	6.3	6.3	5.9
	<i>af</i>	6.0	4.8	4.8	4.7	4.7	4.7	4.6
	<i>au</i>	5.8	5.1	5.0	4.6	4.5	4.5	4.3
	<i>in</i>	6.6	6.0	5.8	5.8	5.8	5.8	5.7
	<i>sg</i>	6.5	5.1	5.1	5.1	5.1	5.1	4.9

- Almost no more improvements after 3 iterations

Conclusion & Discussion

Easy-to-use:

- A pretrained Model + 1-hour ***unlabeled*** speech
- **13.5%** relative WER reduction across **14** target domains (noise, accent, etc.)

Generalization:

- Other models: SeamlessM4T, OWSM, Canary
- Other task: Speech Translation

Anti-forgetting:

- Avoid common catastrophic forgetting in domain adaptation

Discussion

- Large models' attention matrix can present their uncertainty
- Self-improvement is possible in large speech foundation Model

Enhancing Zero-shot Text-to-Speech Synthesis with Human Feedback

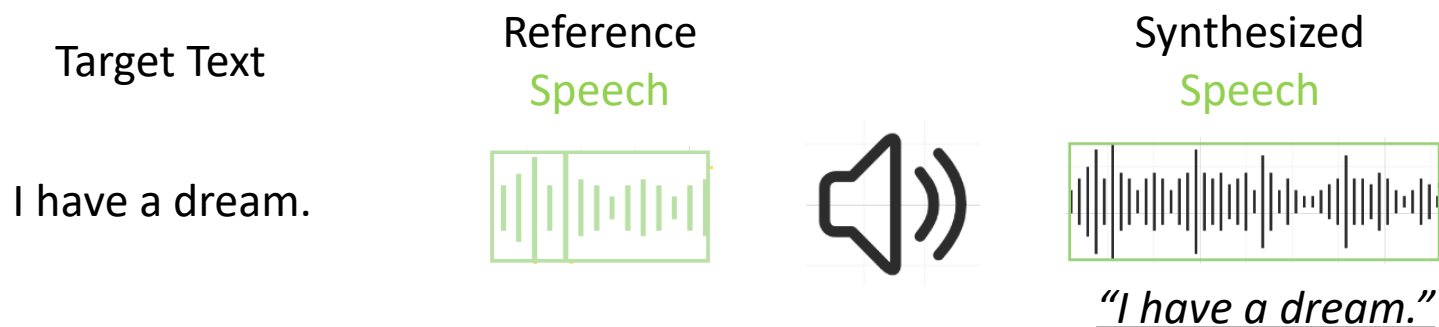
Chen Chen^{1*}, Yuchen Hu^{1*}, Wen Wu², Helin Wang³,
Eng Siong Chng¹, Chao Zhang⁴

¹ Nanyang Technological University ² University of Cambridge
³ Johns Hopkins University ⁴ Tsinghua University



Zero-shot TTS

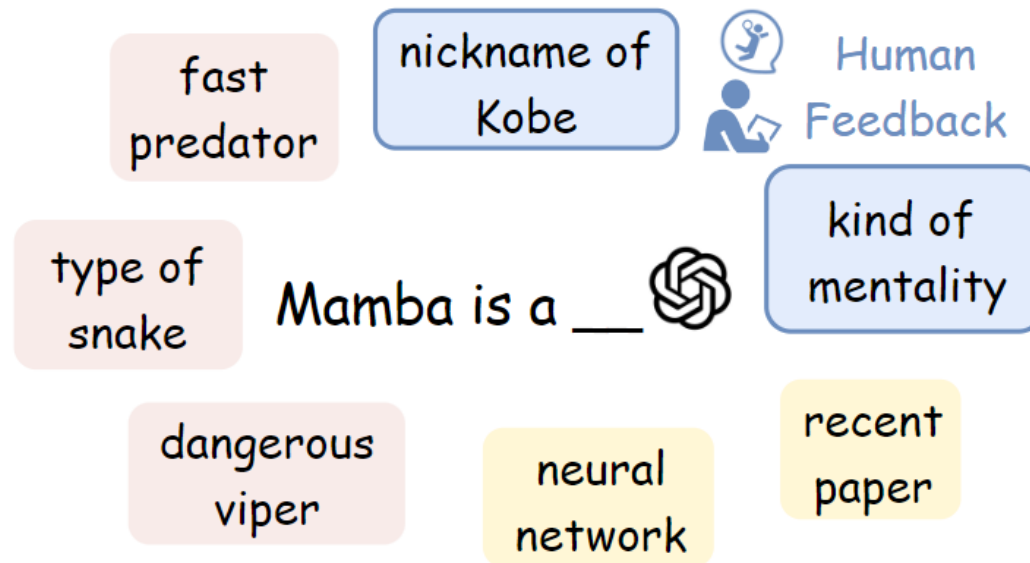
Zero-shot TTS: speak a sentence with *cloned voice*



Current open-sourced TTS model is not robust enough:

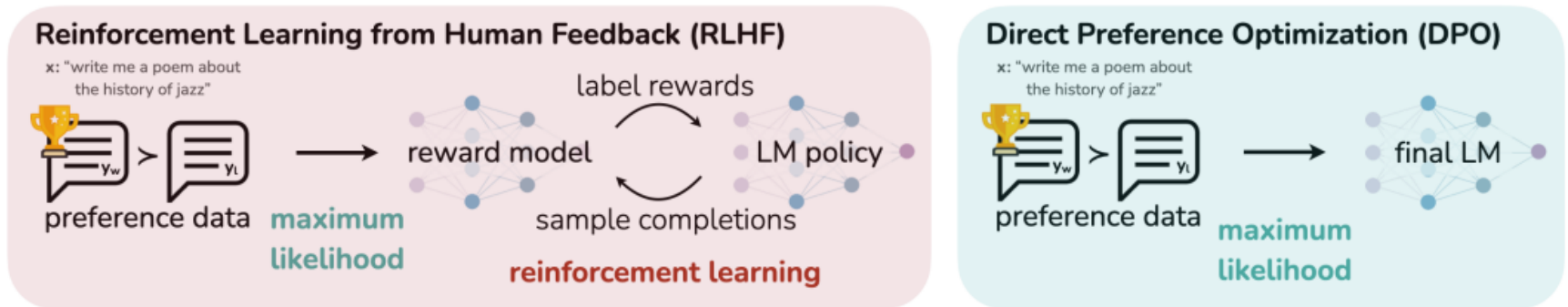
- **Objective Intelligibility:** missing words, wrong words, repetition
- **Subjective Naturalness:** unnatural prosody, tones (human preference)

RLHF is What We Need



- Align LLM's generation with human preference

PPO vs. DPO



- **PPO^[1] (OpenAI)**: most popular RLHF algorithm, make ChatGPT a success
- **DPO^[2] (Stanford)**: direct preference optimization without reward model

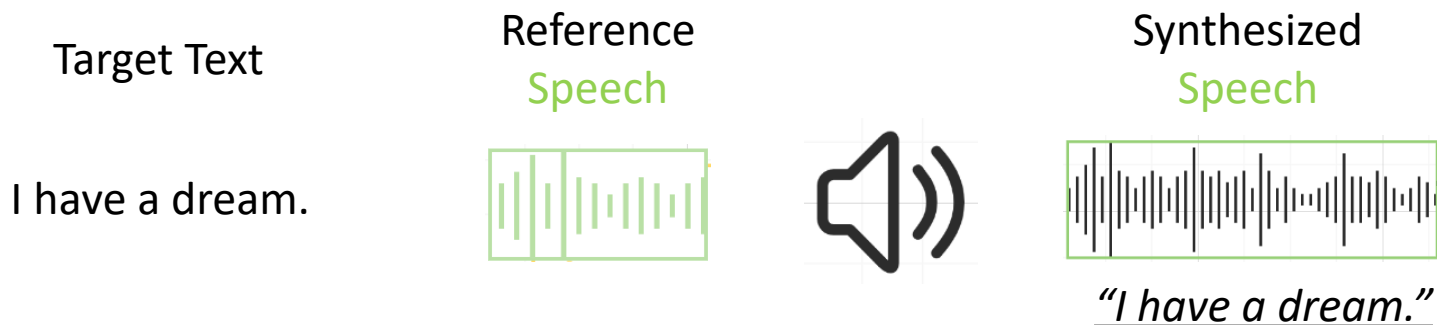
Reward model in TTS: MOSNet --> Not fine-grained enough

[1] Schulman, et al. "Proximal policy optimization algorithms." *arXiv preprint arXiv:1707.06347* (2017).

[2] Rafailov, et al. "Direct preference optimization: Your language model is secretly a reward model." *NeurIPS 2023*.

DPO in TTS?

Zero-shot TTS: speak a sentence with *cloned voice*



Question: If we want to use **DPO**, how to sample binary data?

$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}, \pi_{\text{ref}}) = \mathbb{E} \left[-\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right]$$

Our backbone^[3] lacks diversity:
w and l are not distinctive enough

[3] Peng, et al. "VoiceCraft: Zero-Shot Speech Editing and Text-to-Speech in the Wild." *ACL 2024*.

DPO in TTS?

Question: If we want to use **DPO**, how to sample binary data?

$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}, \pi_{\text{ref}}) =$$

$$\mathbb{E} \left[-\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right]$$

Our backbone^[3] lacks diversity:
w and l are not distinctive enough

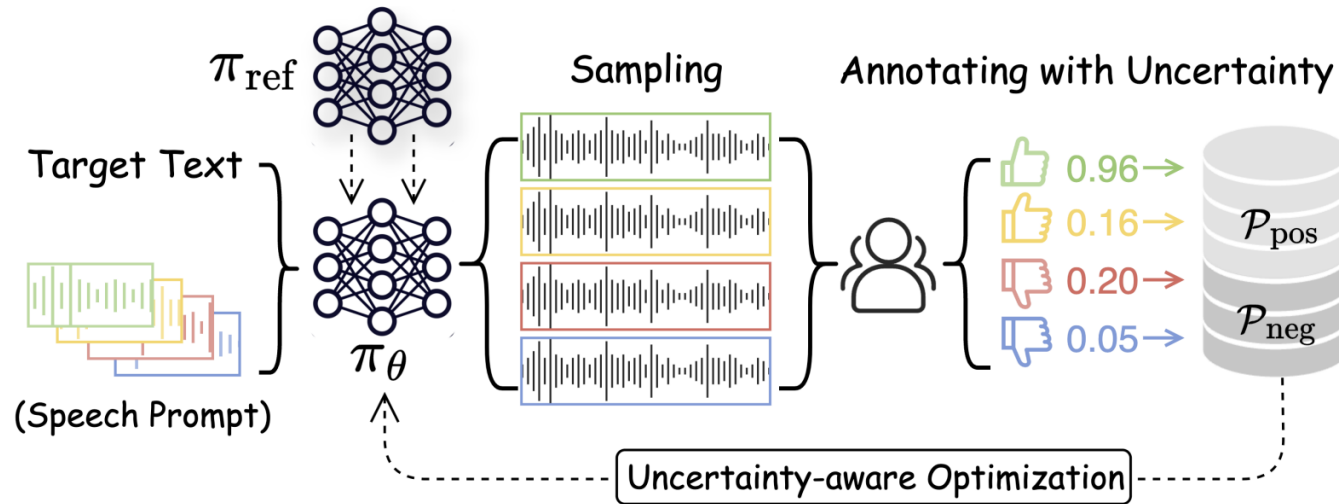
Solutions:

1. Ground-truth as pos, generated as neg (SpeechAlign^[4])
2. Use StyleTTS to produce diversity
3. Change infer hyper-params: top-k, top-p, τ

[4] Zhang, et al. "SpeechAlign: Aligning Speech Generation to Human Preferences." *arXiv preprint arXiv:2404.05600* (2024).

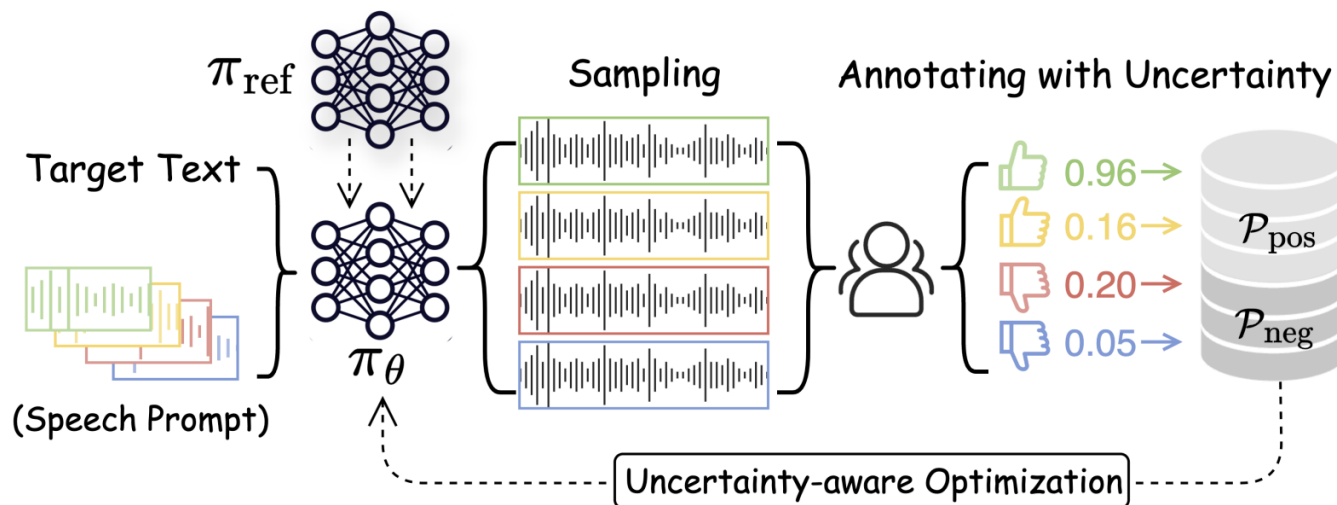


Our Method



- Sampling and annotating with MOS evaluators (human or NN)
- Incorporate MOS's uncertainty to reweight training samples
- Encourage TTS model to synthesize **good** speech and avoid **bad** speech

Our Method



- Encourage TTS model to synthesize **good** speech and avoid **bad** speech

$$\mathcal{L}_{\text{TTS}}(\pi_\theta, \pi_{\text{ref}}) = \mathbb{E}_{t,p,s \sim \mathcal{P}_{\text{pos}} \cup \mathcal{P}_{\text{neg}}} (1 - V_{\text{TTS}}(t, p, s; u)).$$

$$V_{\text{TTS}}(t, p, s; u) = \begin{cases} \sigma(u^{-1} \cdot R(t, p, s) - Z_{\text{ref}}), & \text{if } (t, p, s; u) \sim \mathcal{P}_{\text{pos}} \\ \sigma(Z_{\text{ref}} - u^{-1} \cdot R(t, p, s)), & \text{if } (t, p, s; u) \sim \mathcal{P}_{\text{neg}} \end{cases}$$

$$R(t, p, s) = \log \frac{\pi_\theta(s|t, p)}{\pi_{\text{ref}}(s|t, p)}$$

Objective Evaluation

Model	Label	WER↓ (%)	SIM↑ (0,1)	MOS ↑ by		
				I-CNF	EDL	MOSNet
<i>VoiceCraft (baseline)</i>	-	8.4	0.84	3.51	3.55	3.65
<i>SpeechAlign-DPO</i>	✓	7.2	0.91	3.70 _{+0.19}	3.72 _{+0.17}	3.86 _{+0.21}
<i>SpeechAlign-ODPO</i>	✓	6.9	0.90	3.73 _{+0.21}	3.76 _{+0.24}	3.90 _{+0.25}
<i>PPO-SDP</i>	✗	7.7	0.88	3.65 _{+0.14}	3.69 _{+0.14}	3.85 _{+0.20}
UNO-ICNF	✗	2.6	0.91	3.93_{+0.42}	3.90 _{+0.35}	4.31_{+0.66}
UNO-EDL	✗	2.4	0.92	3.88 _{+0.37}	3.91_{+0.36}	4.28 _{+0.63}
<i>GroundTruth (upper-bound)</i>	-	2.0	-	4.15	4.19	4.52

- **ICNF/EDL**: which model to annotate training data
- **ICNF/EDL**: which model to do evaluation
- **MOSNet**: external evaluation

Training data: 400 utterances (200 pos, 200 neg)

Subjective Evaluation

Table 2: Results on human evaluation.

Model	MOS by	
	Human	MOSNet
<i>VoiceCraft</i>	3.38	3.57
<i>UNO-ICNF</i>	4.06	4.20
<i>UNO-Human</i>	3.98	4.13
<i>GroundTruth</i>	4.55	4.46

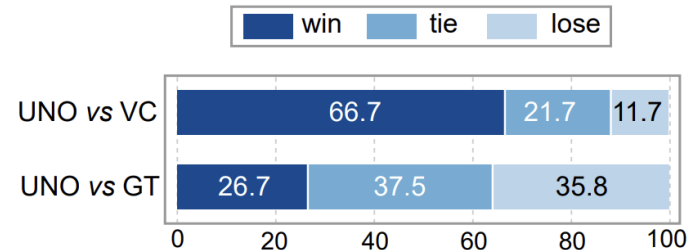



Figure 3: Result of A/B test. “VC” and “GT” denote the “VoiceCraft” and “GroundTruth”.

- **Human**: human annotated training data
- **Human**: human evaluation

Demos


- **Objective Intelligibility:** missing words

VoiceCraft: 

Ours: 

“If a layman in giving baptism pour the water before saying the words is the child baptized”

- **Subjective Naturalness:** unnatural rhythm and tones

VoiceCraft: 

Ours: 

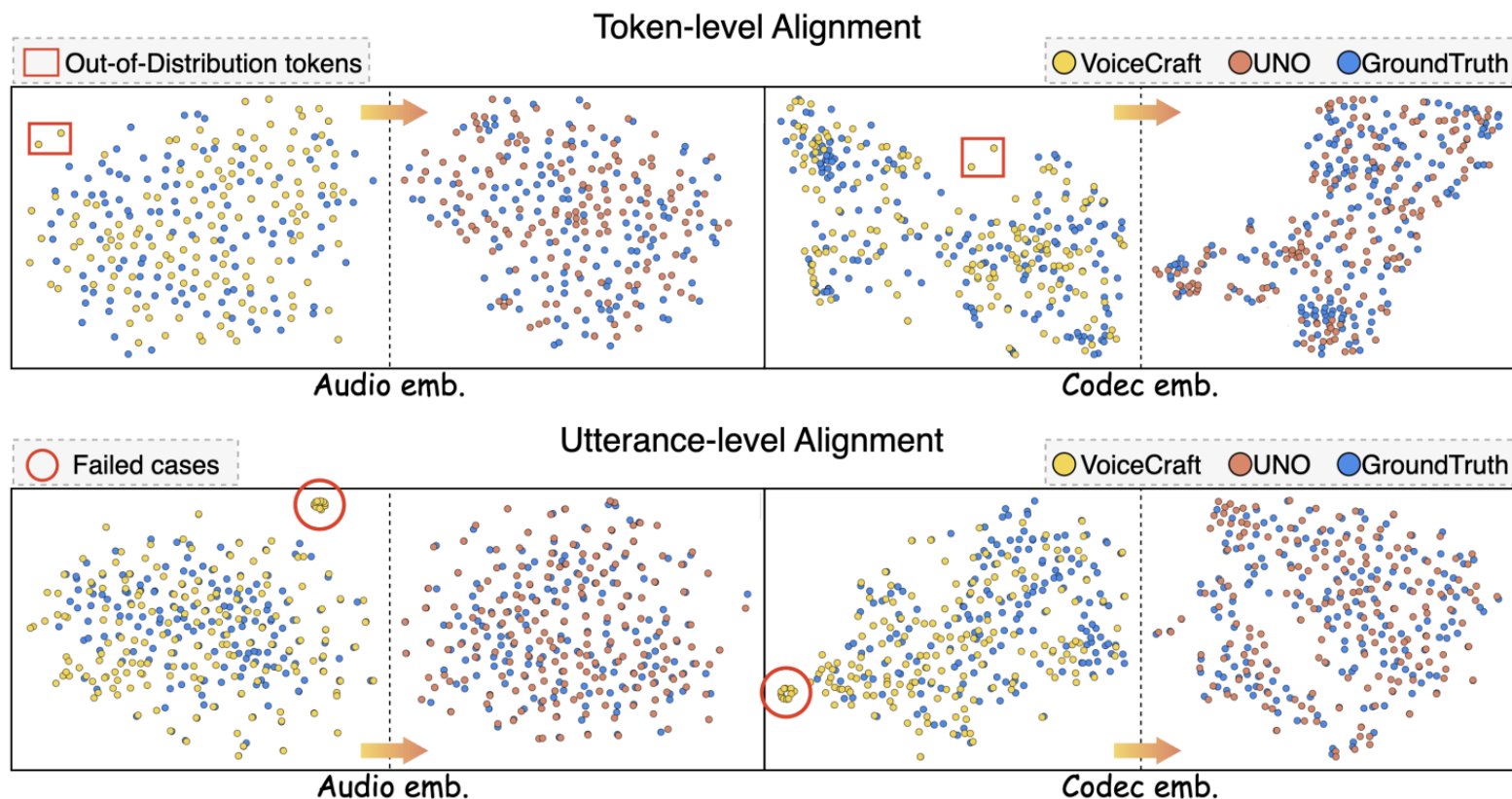
“Not only this but on the table I found a small ball of black dough or clay with specks of something which looks like sawdust in it”

<https://yuchen005.github.io/UNO-TTS-demos/>

28



Visualization



- Align better to the ground-truth speech
- Remove some bad cases

Scalability

EmotionTTS-Valence						EmotionTTS-Arousal					
\mathcal{P}_{pos}		\mathcal{P}_{neg}		$\text{Valence} \uparrow$		\mathcal{P}_{pos}		\mathcal{P}_{neg}		$\text{Arousal} \uparrow$	
\bar{v}	\bar{m}	\bar{v}	\bar{m}	<i>before</i>	<i>after</i>	\bar{a}	\bar{m}	\bar{a}	\bar{m}	<i>before</i>	<i>after</i>
0.65	4.08	0.36	4.04	0.55	0.67	0.69	4.05	0.48	4.20	0.62	0.71

- We can also use RLHF to improve emotion TTS

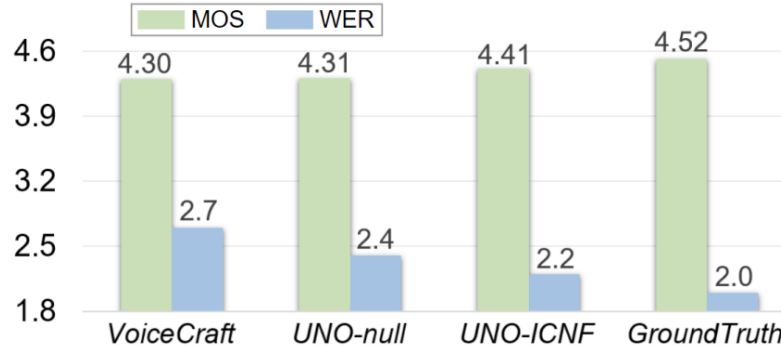
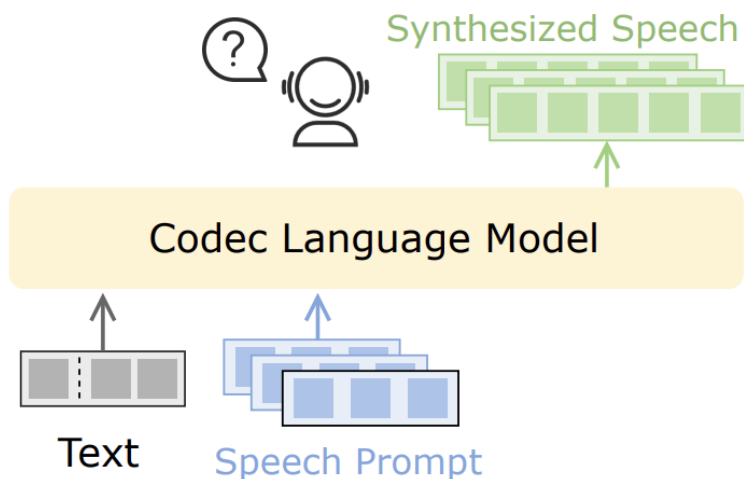


Figure 4: WER and MOS Results on 830M models.

- *VoiceCraft-830M* is already very good
- Our improvements are limited

RLHF on Larger TTS Model

VoiceCraft-830M synthesizes abundant **good** and **non-distinctive** speech:

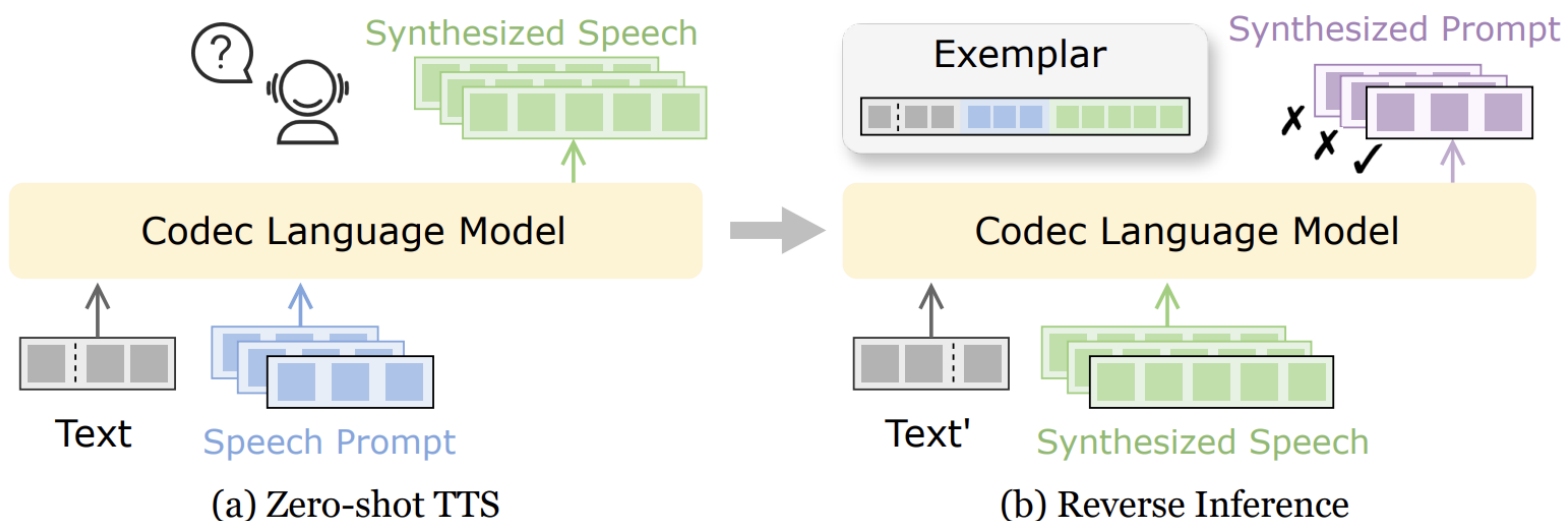


(a) Zero-shot TTS

- Are they really *good* ?
- Does TTS model understand what he is speaking ?

Reverse Inference Optimization

Let TTS model itself perceives the synthesized speech:



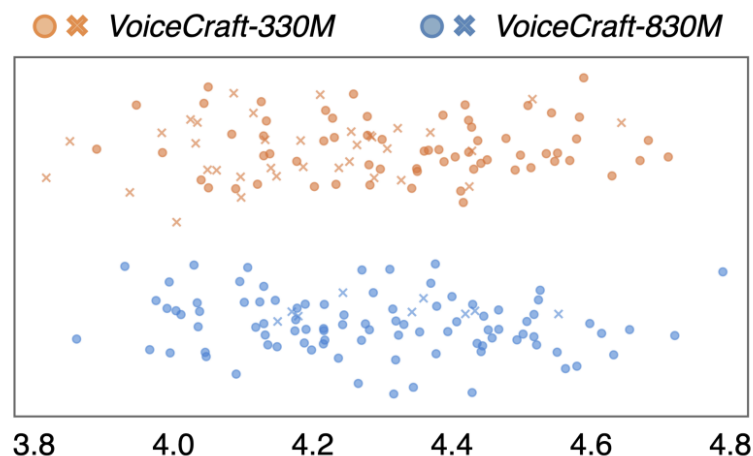
Bayes' formula:

$$P(\mathbf{Y}|\mathbf{T}_Y, \mathbf{T}_X, \mathbf{X}) = \frac{P(\mathbf{X}|\mathbf{T}_X, \mathbf{T}_Y, \mathbf{Y}) P(\mathbf{Y}|\mathbf{T}_Y, \mathbf{T}_X)}{P(\mathbf{X}|\mathbf{T}_X, \mathbf{T}_Y)}$$

$$= \frac{P(\mathbf{Y}|\mathbf{T}_Y)}{P(\mathbf{X}|\mathbf{T}_X)} P(\mathbf{X}|\mathbf{T}_X, \mathbf{T}_Y, \mathbf{Y})$$

Empirical Observation

Ratio of bad reverse inference



MOS vs. Bad reverse inference

Model	✓✗	✓✓
VoiceCraft-330M	3.79	3.90
VoiceCraft-830M	4.24	4.32

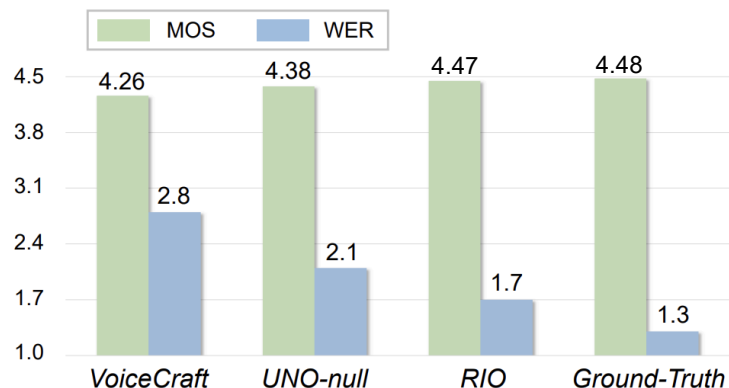
(✓ : good generation, MOS > 3)

- Bad reverse inference results correspond to low TTS robustness
- Human/MOSNet cannot perceive subtle differences in speech, but TTS model can

Evaluation on Longer Test Samples

Model	WER↓ (%)	SIM↑ (0,1)	MOS ↑ by		Bad Case Ratio ↓	
			MOSNet	Human	MOS ≤ 3	%WER > 20
VoiceCraft	35.3	0.79	3.36	3.22	27%	51%
RIO-DPO	11.3	0.92	4.11 _{+0.75}	-	5%	17%
RIO-ODPO	9.2	0.93	4.15 _{+0.79}	-	5%	15%
UNO-null	6.8	0.93	4.20 _{+0.84}	-	4%	11%
RIO (ours)	3.4	0.96	4.40_{+1.04}	4.18_{+0.96}	1%	4%
Ground-Truth	1.3	-	4.48	4.54	0%	0%

Results on 330M models



- Less bad cases, more robust TTS
- Approach GT on 830M backbone

Figure 4: MOS and WER Results on 830M models.

Performance on Large TTS Model

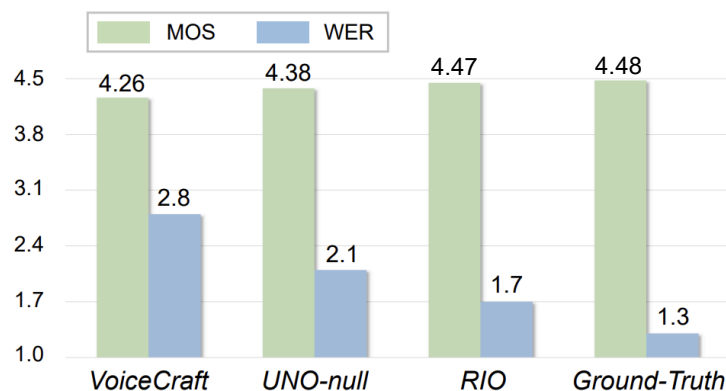


Figure 4: MOS and WER Results on 830M models.

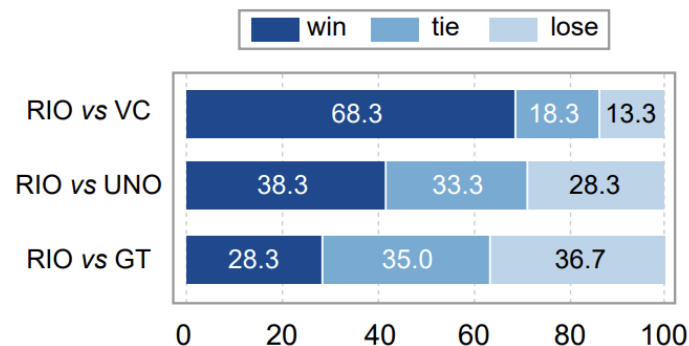
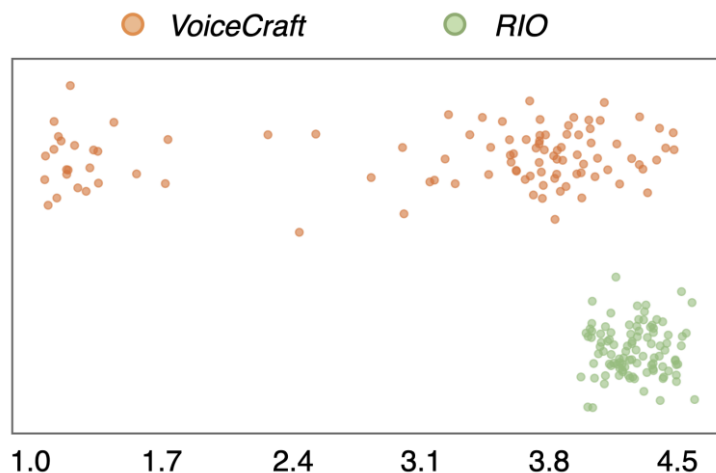


Figure 3: Results of A/B test. “VC” and “GT” denote the “VoiceCraft” and “Ground-Truth”.

- High-level sampling can push the limit of RLHF on large TTS models

Analysis

MOS distribution



Ratio of good reverse inference
in good zero-shot samples

Model	Baseline	RIO
<i>VoiceCraft-330M</i>	54%	85%
<i>VoiceCraft-830M</i>	80%	97%

- Better MOS distribution: higher mean and lower variance
- More good reverse inference results after post-training --> better robustness

Conclusion & Discussion

RLHF Post-training in TTS

- Improve subjective naturalness and expressiveness
- Correct many objective issues

Annotating is important in RLHF-TTS:

- Positive samples: Good enough, i.e., natural, expressive, correct
- Negative samples: Cover as many error types as possible

Discussion:

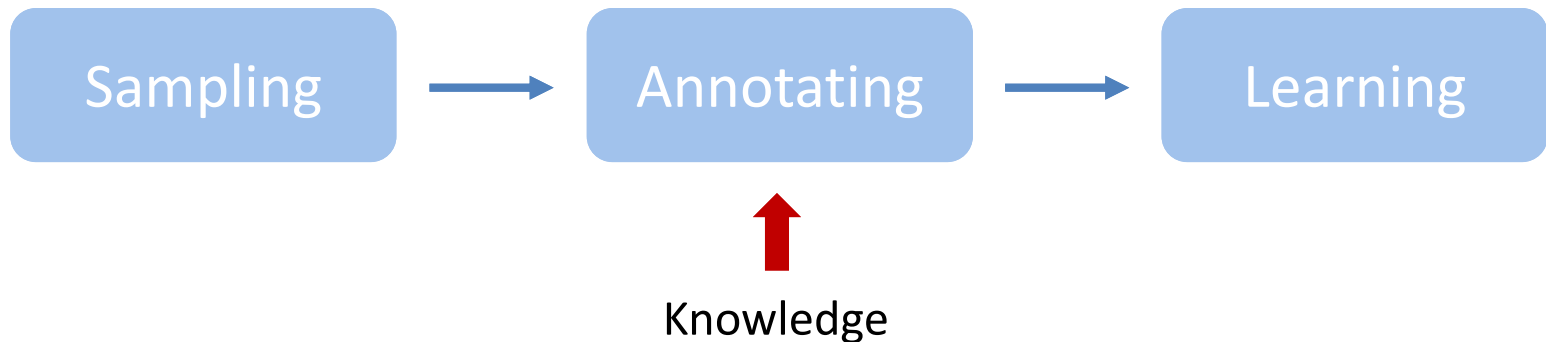
- For general speech perception: Model > Human
- TTS model can **not** understand what it predict, similar to Whisper/Qwen-Audio
- Future: Speech/text understanding/generation are unified in one single model

Summary

Pretraining & SFT:

$$\mathcal{L}(x, y) = \sum_{l=1}^L -\log \mathcal{P}_{\theta}(y_l | y_{l-1}, \dots, y_1, x) \quad \text{Learn a distribution --> uncertain}$$

Post-training:



- Appropriate annotating can help efficiently improve the performance!

Thank you! & QA



Appendix: LLM Hallucination



Please describe this image in detail.

Human

Beam Search

LLaVA-1.5

The image features a blue bowl filled with a delicious mixture of bananas, nuts, and oatmeal. The bowl is placed on a dining table, and a spoon is resting inside the bowl, ready to be used for enjoying the meal. In addition to the bowl of food, there are a few other items on the table. A bottle can be seen on the left side of the table, while a cup is positioned towards the top right corner. A book is also present on the right side of the table, adding to the cozy atmosphere of the scene.

Hallucination!

Our OPERA

LLaVA-1.5

The image features a blue bowl filled with a delicious mixture of bananas, nuts, and oatmeal. The bowl is placed on a dining table, and a spoon is resting inside the bowl, ready to be used. The bananas are scattered throughout the bowl, with some closer to the top and others near the bottom. The nuts and oatmeal complement the bananas, creating a visually appealing and appetizing dish.

- **Non-Hallucinations:** describes the food (e.g., bananas, nuts, oatmeal) inside the bowl
- **Hallucinations:** imagines the items on the table that is outside the image

NOTE: Hallucinations starts with “In addition to ...”

Appendix: LLM Hallucination

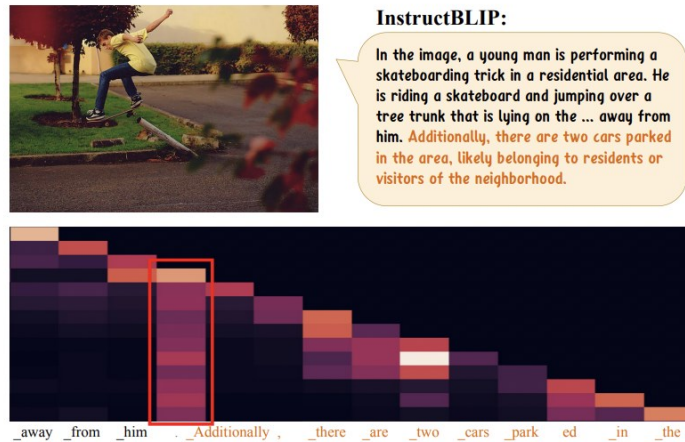


Figure 2. A case of relationship between hallucinations and knowledge aggregation patterns. Hallucinations are highlighted.

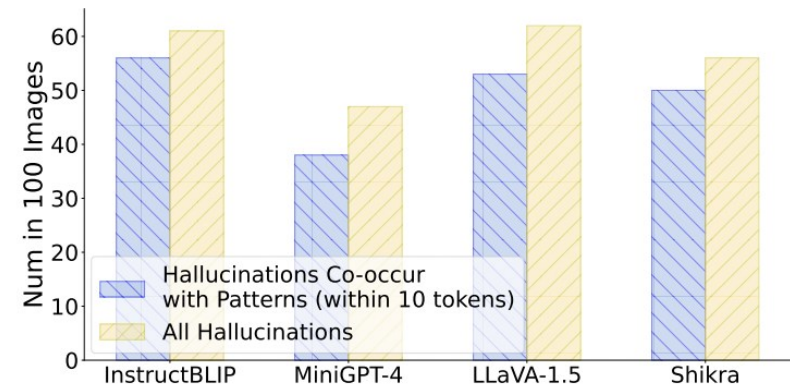


Figure 3. Hallucinations often start within the first 10 tokens after knowledge aggregation patterns.

- Hallucinations are usually triggered by specific tokens (e.g., “*additionally*”);
- We can observe a “knowledge aggregation pattern” in self-attention map along with the beginning of hallucinations → *An insightful finding!*

Appendix: LLM Hallucination



InstructBLIP:

In the image, a young man is performing a skateboarding trick in a residential area. He is riding a skateboard and jumping over a tree trunk that is lying on the ... away from him. Additionally, there are two cars parked in the area, likely belonging to residents or visitors of the neighborhood.

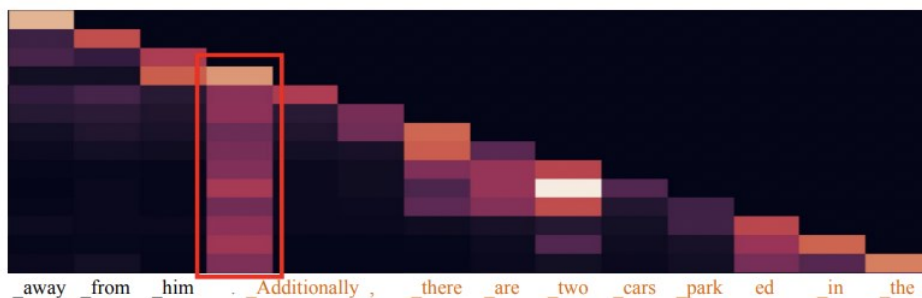


Figure 2. A case of relationship between hallucinations and knowledge aggregation patterns. Hallucinations are highlighted.

All hallucinations are highly related to the starting token “*additionally*” but unrelated to previous normal tokens!