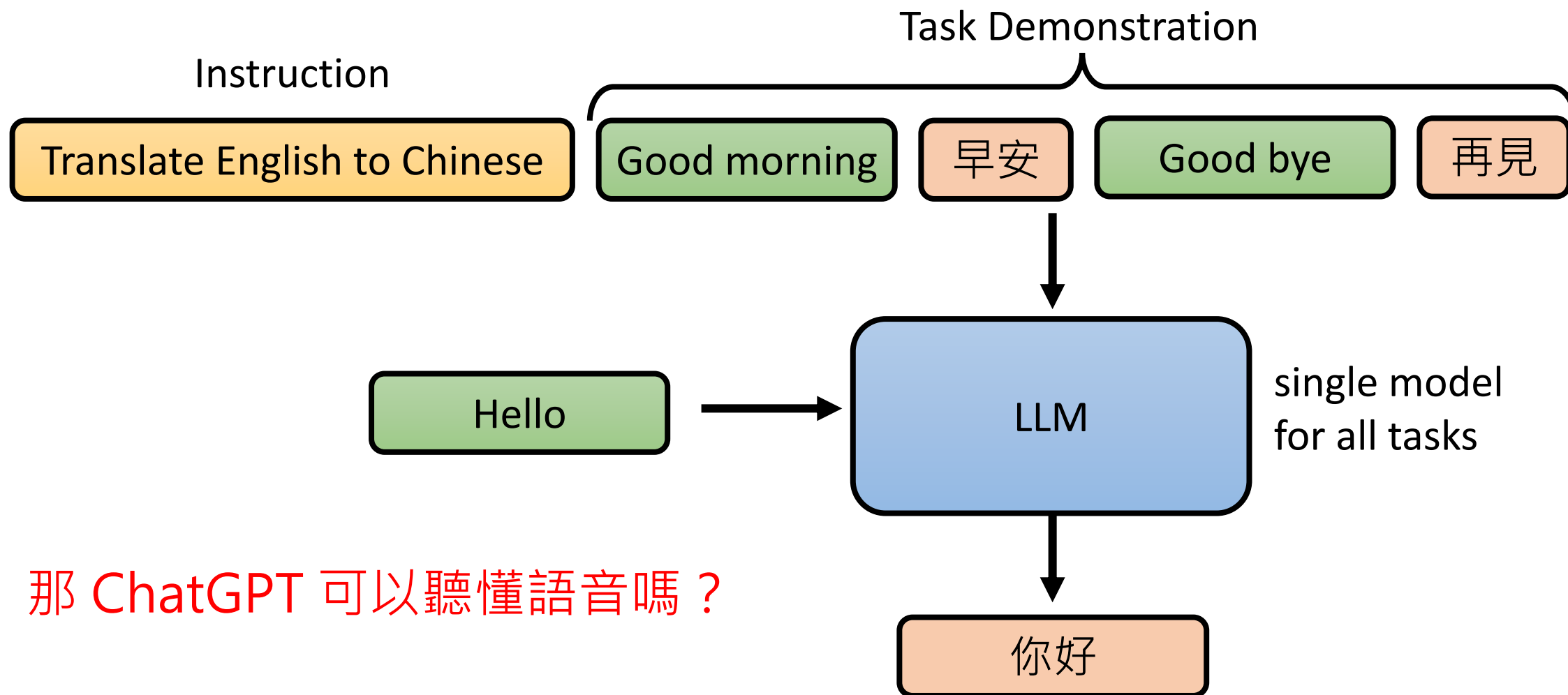


A light blue background featuring several faint, overlapping speech bubbles of various shapes and sizes. A thin, dark diagonal line is positioned in the upper left corner.

Towards a Speech Version of ChatGPT

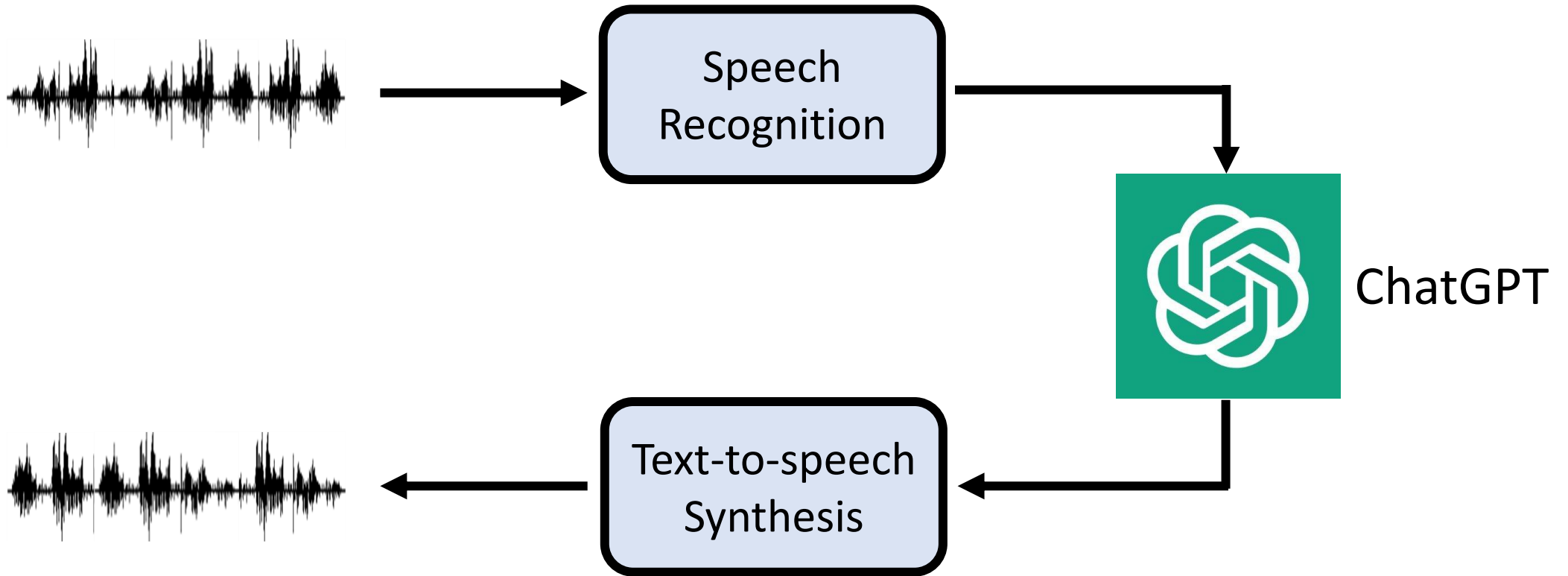
Hung-yi Lee

大家都已經見識到 ChatGPT 的能力



那 ChatGPT 可以聽懂語音嗎？

語音版 ChatGPT? 這不是已經有了嗎？



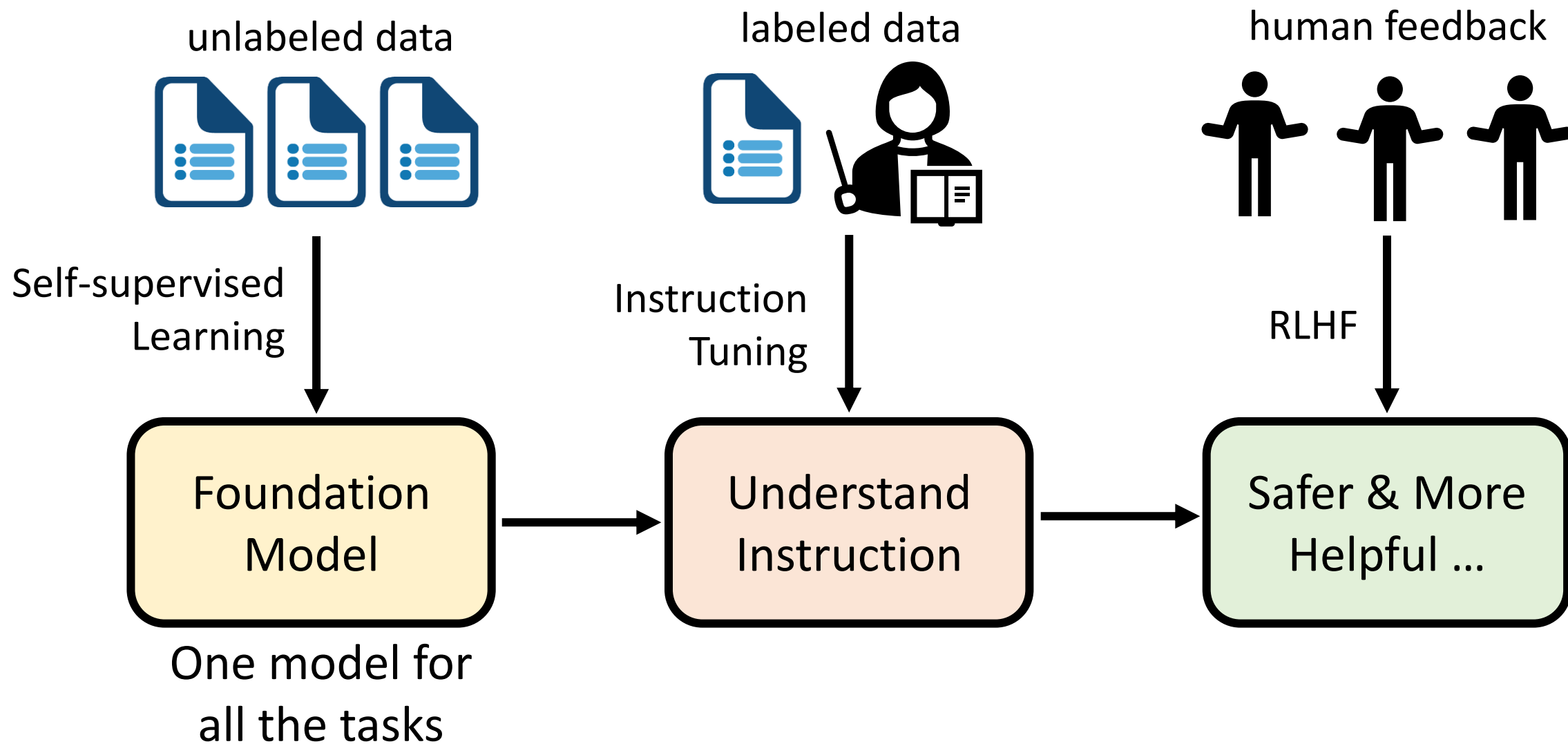
At the end of this talk, I will show you that this trivial solution is insufficient.

我們離語音版的 ChatGPT 其實還有距離



(獵人第八卷)

How did LLM get there?

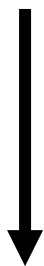


How about Speech?

unlabeled data



Self-supervised
Learning

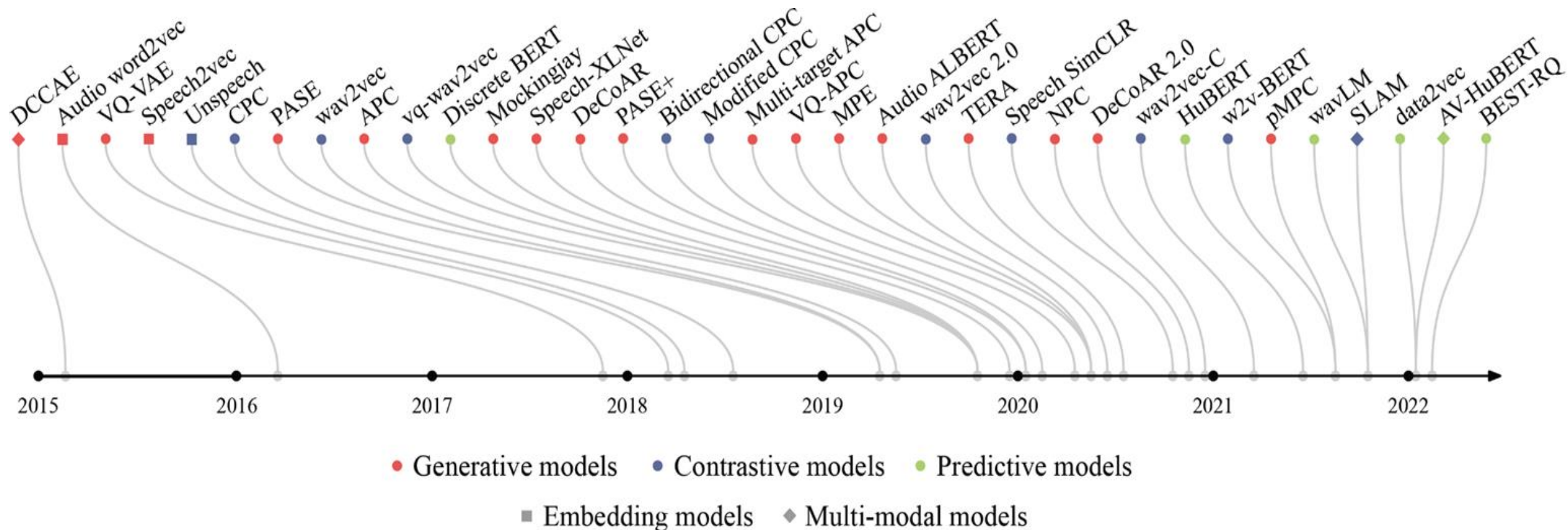


Speech
Foundation
Model

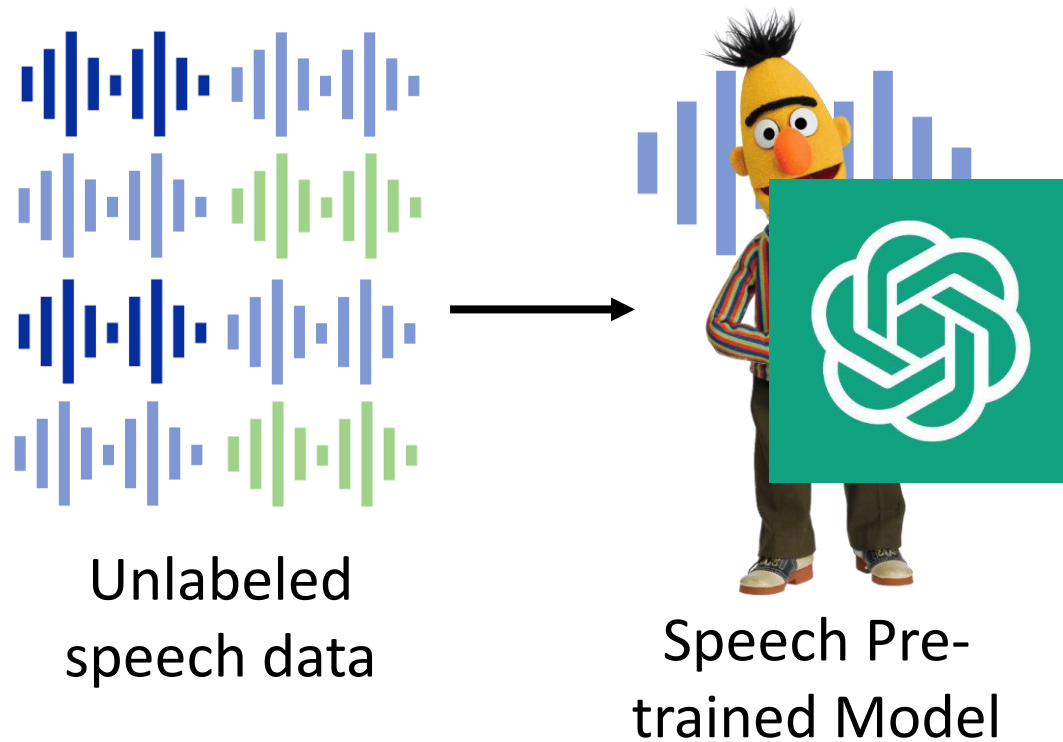
Self-Supervised Speech Representation Learning: A Review

Abdelrahman Mohamed*, Hung-yi Lee*, Lasse Borgholt*, Jakob D. Havtorn*, Joakim Edin, Christian Igel
Katrin Kirchhoff, Shang-Wen Li, Karen Livescu, Lars Maaløe, Tara N. Sainath, Shinji Watanabe

<https://arxiv.org/abs/2205.10643>



Self-supervised Learning for Speech



- There are so many self-supervised models
- They have shown to achieve good performance on ASR.
- Are they specialist for ASR? Or are they universal?

I believe they are specialist.



My two cents
(before 2021)



SUPERB

Speech processing Universal PERformance Benchmark

SUPERB: Speech processing Universal PERformance Benchmark

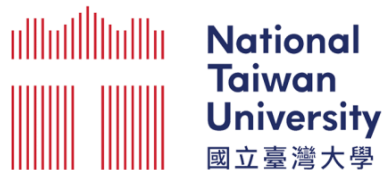


Shu-wen
Yang

Shu-wen Yang¹, Po-Han Chi^{1}, Yung-Sung Chuang^{1*}, Cheng-I Jeff Lai^{2*}, Kushal Lakhota^{3*},
Yist Y. Lin^{1*}, Andy T. Liu^{1*}, Jiatong Shi^{4*}, Xuankai Chang⁶, Guan-Ting Lin¹,
Tzu-Hsien Huang¹, Wei-Cheng Tseng¹, Ko-tik Lee¹, Da-Rong Liu¹, Zili Huang⁴, Shuyan Dong^{5†},
Shang-Wen Li^{5†}, Shinji Watanabe⁶, Abdelrahman Mohamed³, Hung-yi Lee¹*

<https://arxiv.org/abs/2105.01051>

INTERSPEECH 2021



JOHNS HOPKINS
UNIVERSITY

Carnegie
Mellon
University



Massachusetts
Institute of
Technology

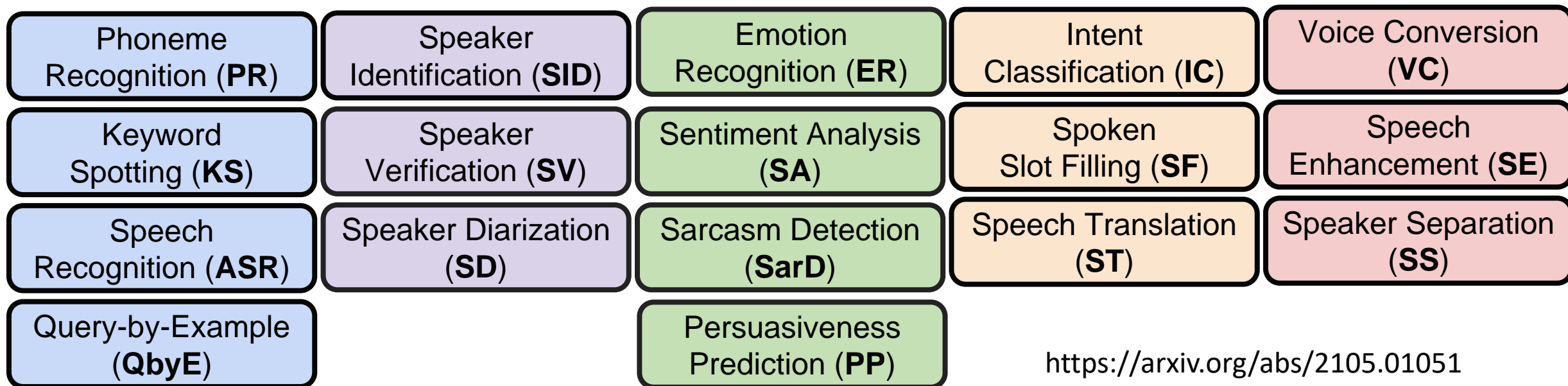




SUPERB

Speech processing Universal PERformance Benchmark

Evaluate a wide range of speech self-supervised models on many speech tasks



Content



Speaker



Prosody



Semantic



Synthesis

<https://arxiv.org/abs/2105.01051>

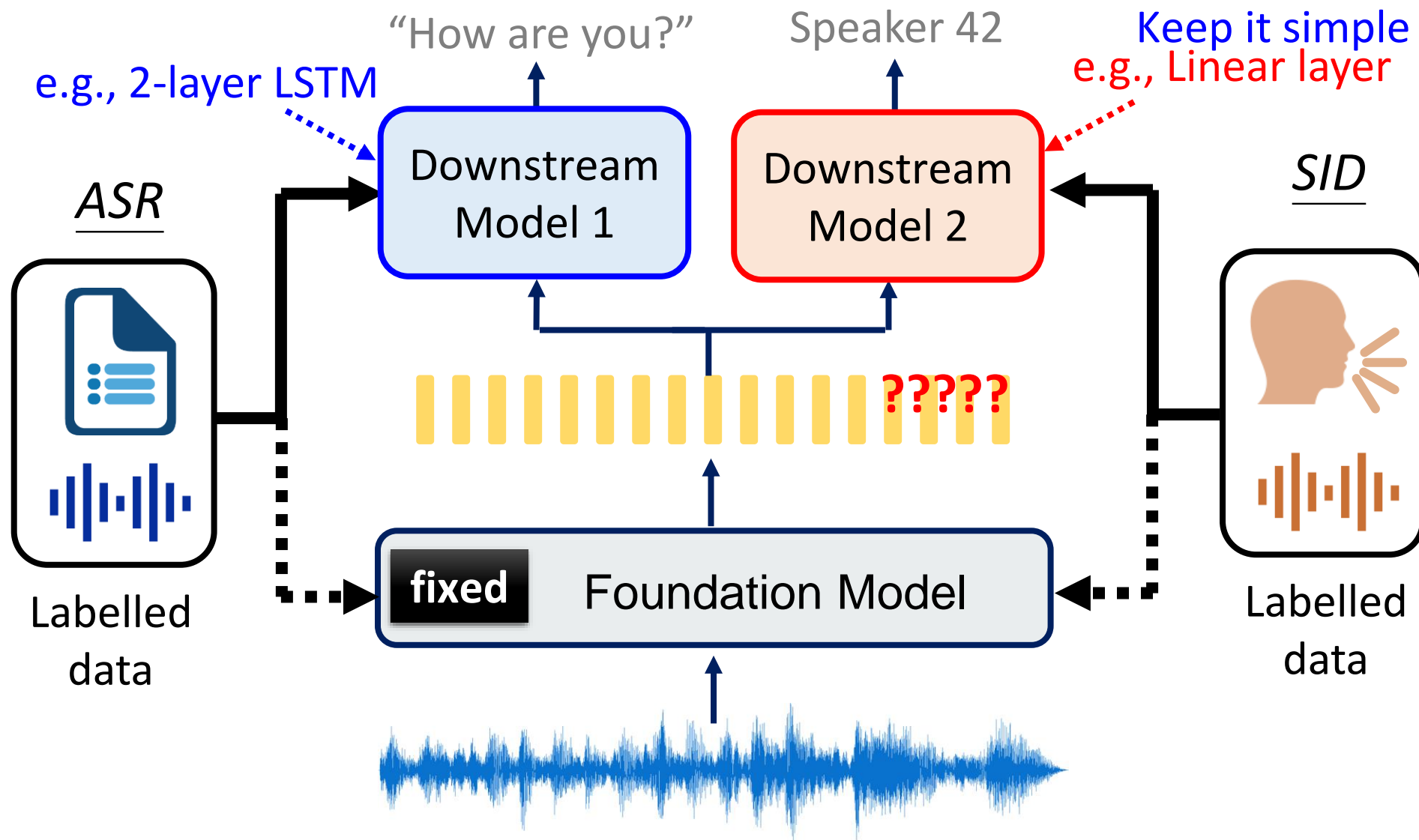
<https://arxiv.org/abs/2203.06849>

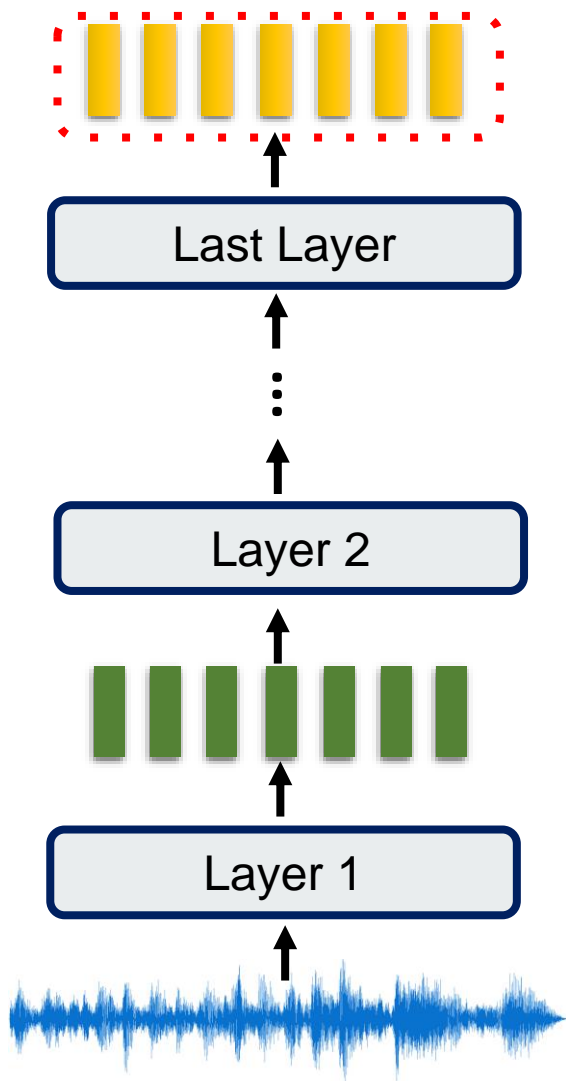
<https://arxiv.org/abs/2210.08634>

<https://arxiv.org/abs/2210.07185>

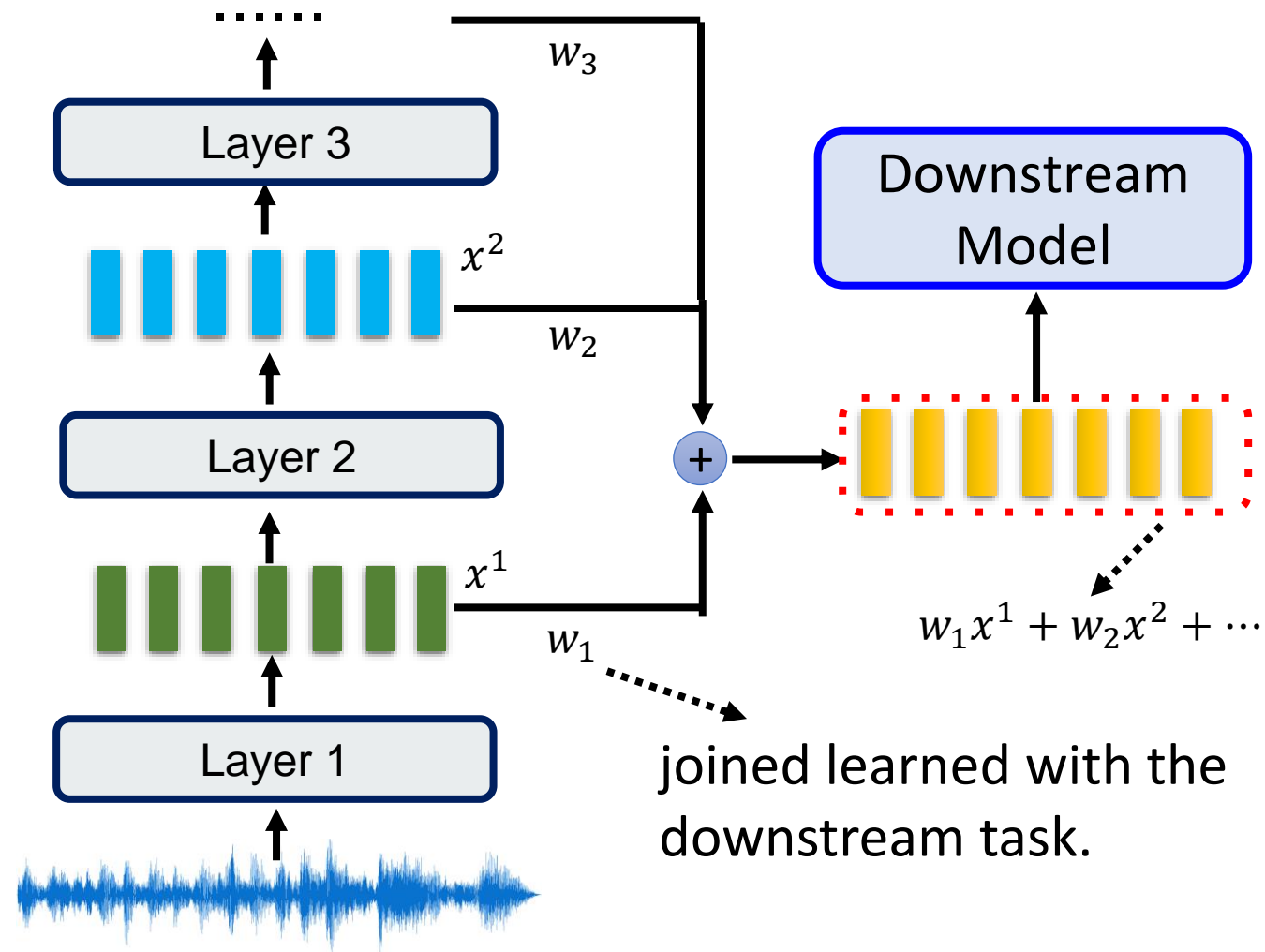
<https://arxiv.org/abs/2110.06280>

How to use Speech Foundation Model





Not always lead to decent performance.



joined learned with the downstream task.

“Weighted-sum” is Very effective!

Results of SUPERB

	Content				Speaker			Semantic		Emotion
	PR	KS	ASR	QbE	SID	ASV	SD	IC	SF	ER
fbank	82.01	8.63	15.21	0.0058	8.50E-04	9.56	10.05	9.1	69.64	35.39
PASE+	58.87	82.54	16.62	0.0072	37.99	11.61	8.68	29.82	62.14	57.86
APC	41.98	91.01	14.74	0.0310	60.42	8.56	10.53	74.69	70.46	59.33
VQ-APC	41.08	91.11	15.21	0.0251	60.15	8.72	10.45	74.48	68.53	59.66
NPC	43.81	88.96	13.91	0.0246	55.92	9.40	9.34	69.44	72.79	59.08
Mockingjay	70.19	83.67	15.48	6.60E-04	32.29	11.66	10.54	34.33	61.59	50.28
TERA	49.17	89.48	12.16	0.0013	57.57	15.89	9.96	58.42	67.50	56.27
DeCoAR 2.0	14.93	94.48	9.07	0.0406	74.42	7.16	6.59	90.80	83.28	62.47
modified CPC	42.54	91.88	13.53	0.0326	39.63	12.86	10.38	64.09	71.19	60.96
wav2vec	31.58	95.59	11.00	0.0485	56.56	7.99	9.90	84.92	76.37	59.79
vq-wav2vec	33.48	93.38	12.80	0.0410	38.80	10.38	9.93	85.68	77.68	58.24
wav2vec 2.0 base	5.74	96.23	4.79	0.0233	75.18	6.02	6.08	92.35	88.30	63.43
wav2vec 2.0 large	4.75	96.66	3.10	0.0489	86.14	5.65	5.62	95.28	87.11	65.64
HuBERT base	5.41	96.30	4.79	0.0736	81.42	5.11	5.88	98.34	88.53	64.92
HuBERT large	3.53	95.29	2.94	0.0353	90.33	5.98	5.75	98.76	89.81	67.62

<https://arxiv.org/abs/2105.01051>

Results of SUPERB

	Content				Speaker			Semantic		Emotion
	PR	KS	ASR	QbE	SID	ASV	SD	IC	SF	ER
fbank	82.01	8.63	15.21	0.0058	8.50E-04	9.56	10.05	9.1	69.64	35.39
PASE+	58.87	82.54		0.0072	37.99		8.68	29.82		57.86
APC	41.98	91.01	14.74	0.0310	60.42	8.56		74.69	70.46	59.33
VQ-APC	41.08	91.11		0.0251	60.15	8.72		74.48		59.66
NPC	43.81	88.96	13.91	0.0246	55.92	9.40	9.34	69.44	72.79	59.08
Mockingjay	70.19	83.67			32.29			34.33		50.28
TERA	49.17	89.48	12.16		57.57		9.96	58.42		56.27
DeCoAR 2.0	14.93	94.48	9.07	0.0406	74.42	7.16	6.59	90.80	83.28	62.47
modified CPC	42.54	91.88	13.53	0.0326	39.63			64.09	71.19	60.96
wav2vec	31.58	95.59	11.00	0.0485	56.56	7.99	9.90	84.92	76.37	59.79
vq-wav2vec	33.48	93.38	12.80	0.0410	38.80		9.93	85.68	77.68	58.24
wav2vec 2.0 base	5.74	96.23	4.79	0.0233	75.18	6.02	6.08	92.35	88.30	63.43
wav2vec 2.0 large	4.75	96.66	3.10	0.0489	86.14	5.65	5.62	95.28	87.11	65.64
HuBERT base	5.41	96.30	4.79	0.0736	81.42	5.11	5.88	98.34	88.53	64.92
HuBERT large	3.53	95.29	2.94	0.0353	90.33	5.98				

They can be universal!

To learn more



Compression



Robust



Adapter/Prompt



Unsupervised
ASR

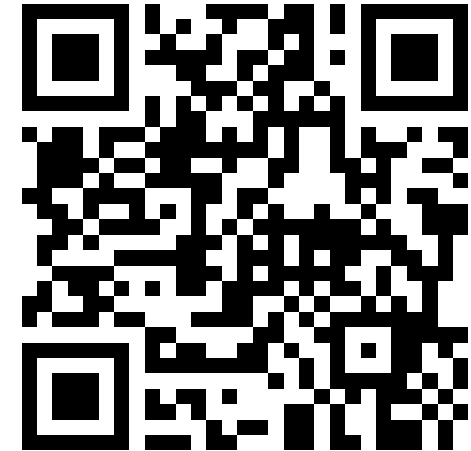


Visual-enhanced



Prosody

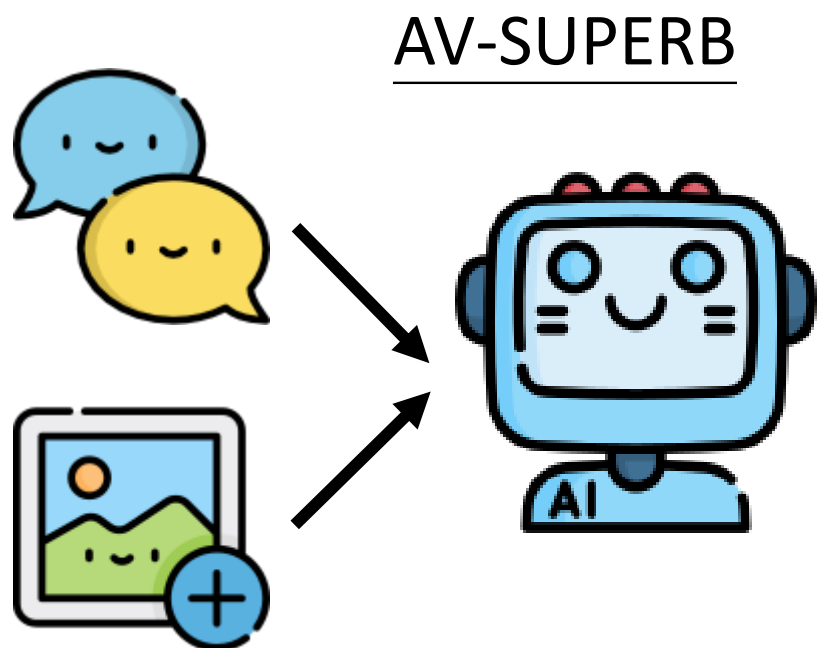
Research Group
@ JSALT 2022



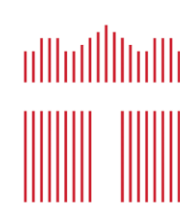
https://youtu.be/_GbZRM18NxQ

Closing-day
Presentation

More SUPERB families are coming



<https://av.superbbenchmark.org/>



National
Taiwan
University
國立臺灣大學

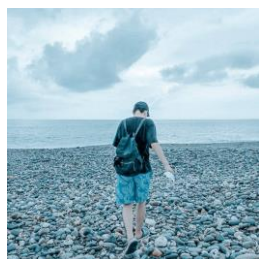


Carnegie
Mellon
University

Meta

REMBRAND

AV-SUPERB: A MULTI-TASK EVALUATION BENCHMARK FOR AUDIO-VISUAL REPRESENTATION MODELS



Yuan Tseng

Yuan Tseng¹, Layne Berry^{2}, Yi-Ting Chen^{3*}, I-Hsiang Chiu^{1*}, Hsuan-Hao Lin^{1*}, Max Liu^{1*},
Puyuan Peng^{2*}, Yi-Jen Shih^{1*}, Hung-Yu Wang^{1*}, Haibin Wu^{1*}, Po-Yao Huang⁴, Chun-Mao Lai¹,
Shang-Wen Li⁴, David Harwath², Yu Tsao³, Shinji Watanabe⁵, Abdelrahman Mohamed⁶,
Chi-Luen Feng¹, Hung-yi Lee¹*

<https://arxiv.org/abs/2309.10787>

Spoken Language Understanding Evaluation (SLUE) v2

<https://arxiv.org/abs/2212.10525>

SLUE Phase-2: A Benchmark Suite of Diverse Spoken Language Understanding Tasks

**Suwon Shon¹, Siddhant Arora^{2*}, Chyi-Jiunn Lin^{*3}, Ankita Pasad^{4*},
Felix Wu¹, Roshan Sharma², Wei-Lun Wu³,
Hung-Yi Lee³, Karen Livescu⁴, Shinji Watanabe²**

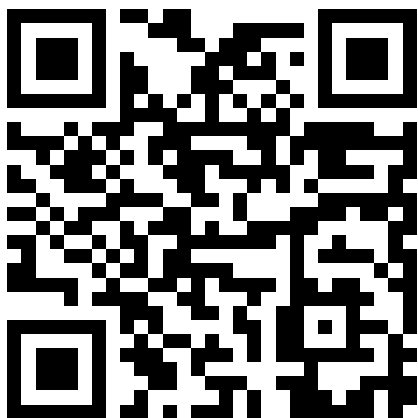
¹ASAPP ²Carnegie Mellon University ³National Taiwan University
⁴Toyota Technological Institute at Chicago

New spoken QA corpus – SLUE-SQA-5: 47k spoken question-context pairs

Real speech!

Useful Toolkit!

The S3PRL toolkit



<https://github.com/s3prl/s3prl>

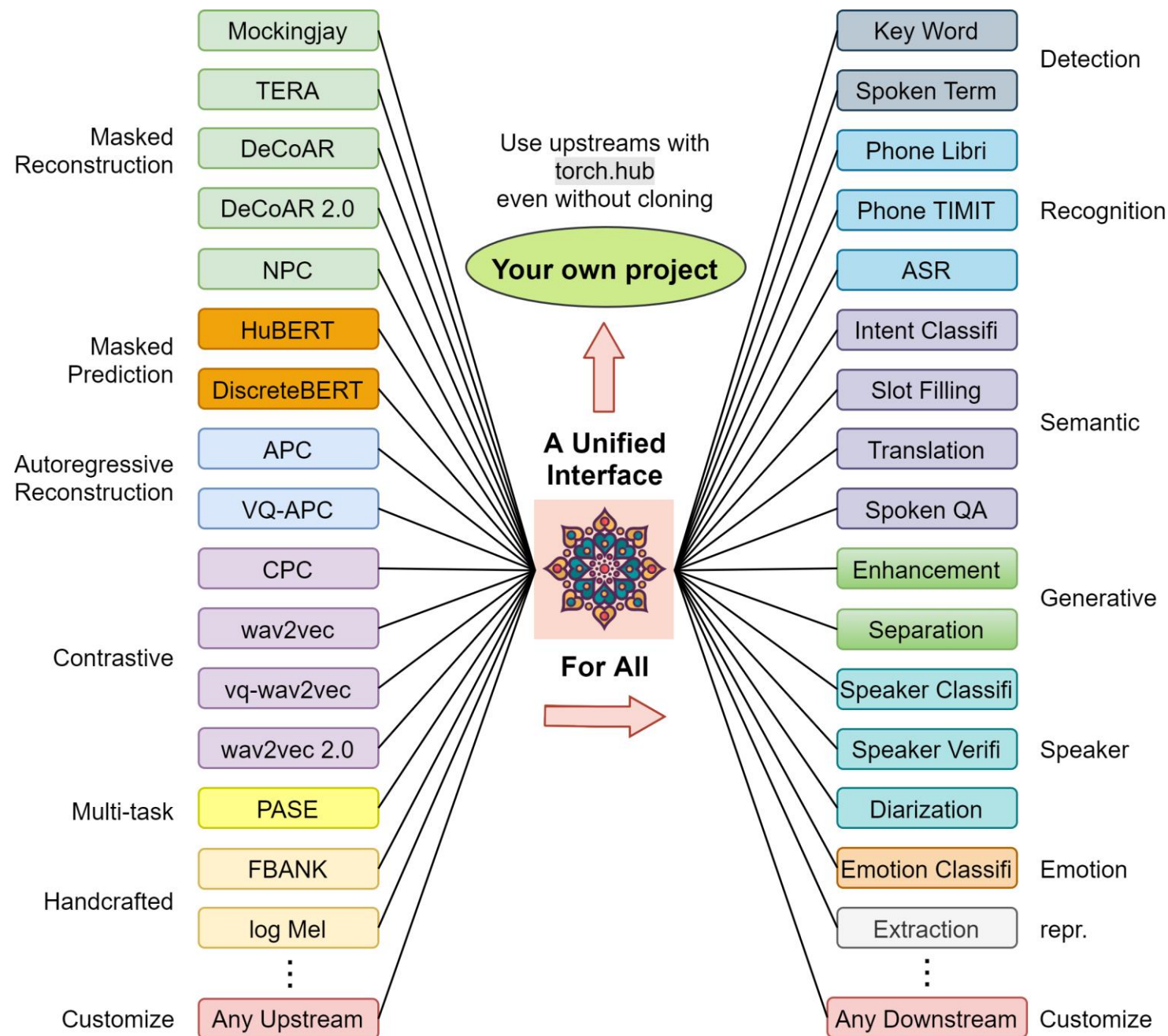
Creator



Shu-wen Yang



Andy T. Liu



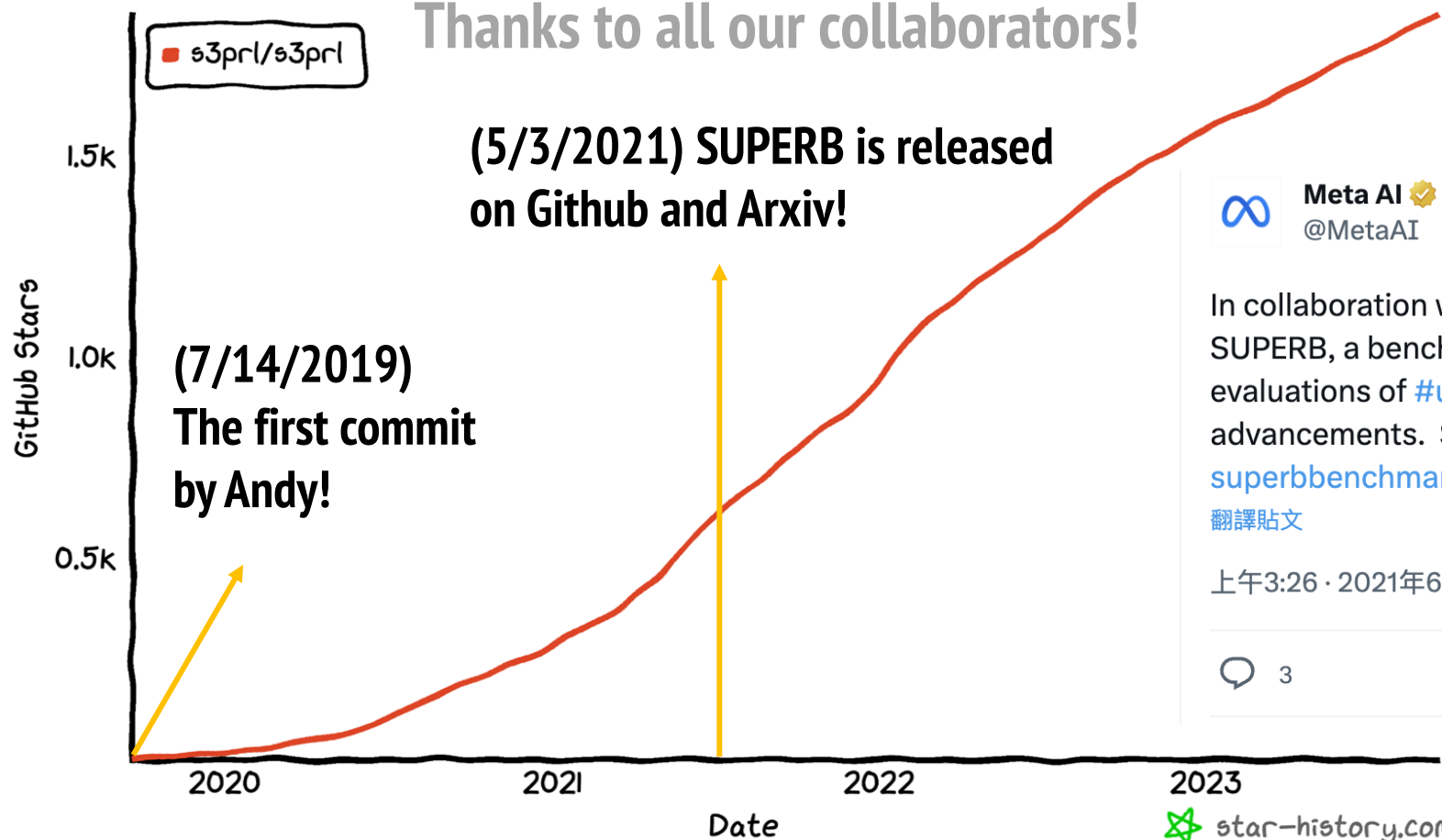
Over 1.9k stars & used by over 48 repos



Shu-wen Yang

Star History

Thanks to all our collaborators!



Meta AI
@MetaAI

(6/19/2021)

Meta helped promote SUPERB!

In collaboration with @ntu_spml, @LTIatCMU, & @jhucisp we introduce SUPERB, a benchmark using 10 speech processing tasks to standardize evaluations of #unsupervised models used in speech processing advancements. Submit & evaluate your models here:

superbenchmark.org

翻譯貼文

上午3:26 · 2021年6月19日



3



60



169



18



star-history.com

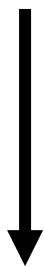
data from
<https://star-history.com/>

How about Speech?

unlabeled data



Self-supervised
Learning



Speech
Foundation
Model

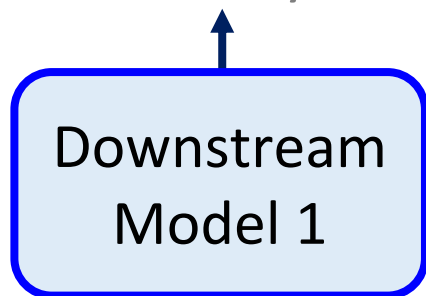
One model for
all the tasks

?

SUPERB

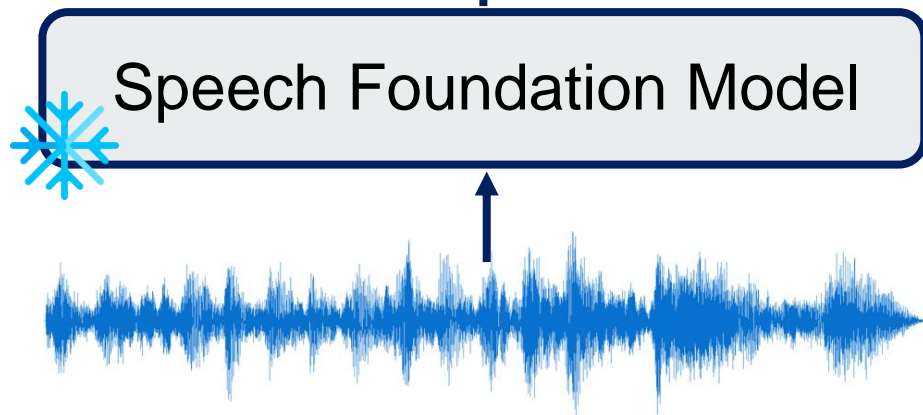
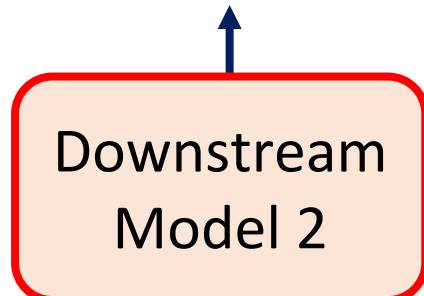
ASR

“How are you?”



SID

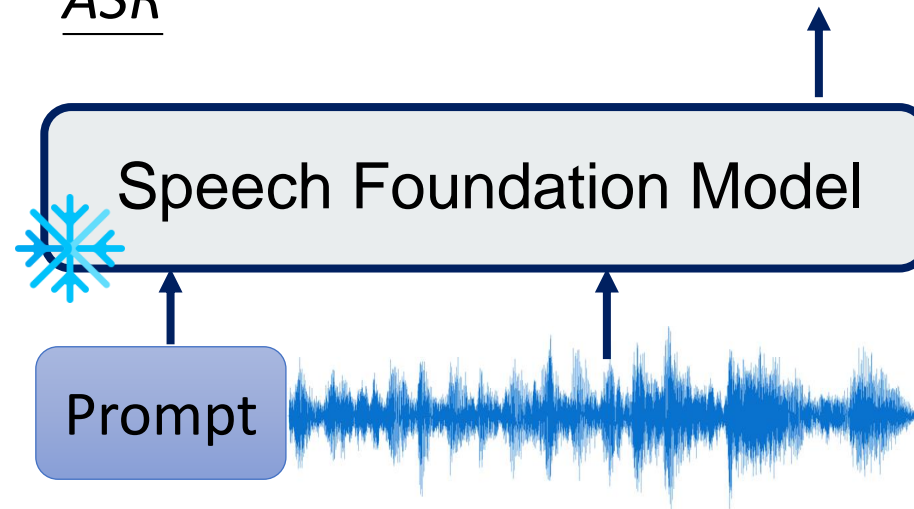
Speaker 42



Prompting

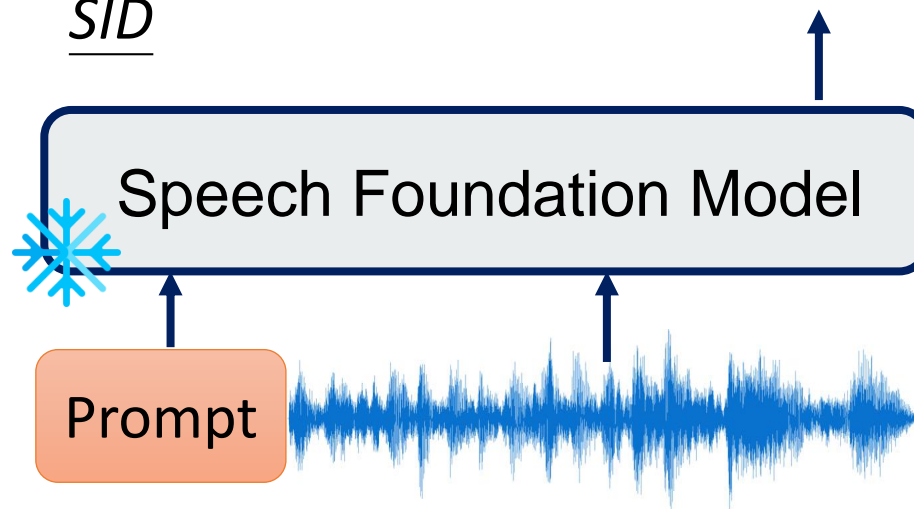
ASR

“How are you?”



SID

Speaker 42

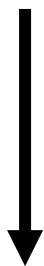


How about Speech?

unlabeled data



Self-supervised
Learning



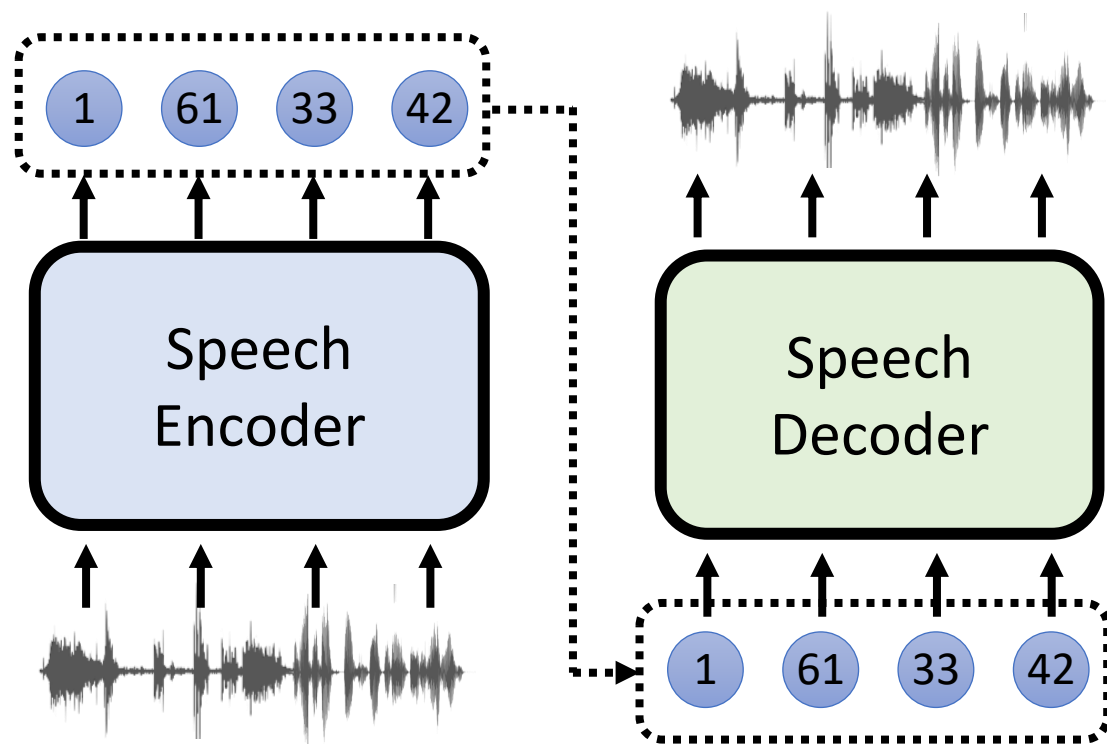
Speech
Foundation
Model

One model for
all the tasks

?

Prompting Speech LM

Speech LM

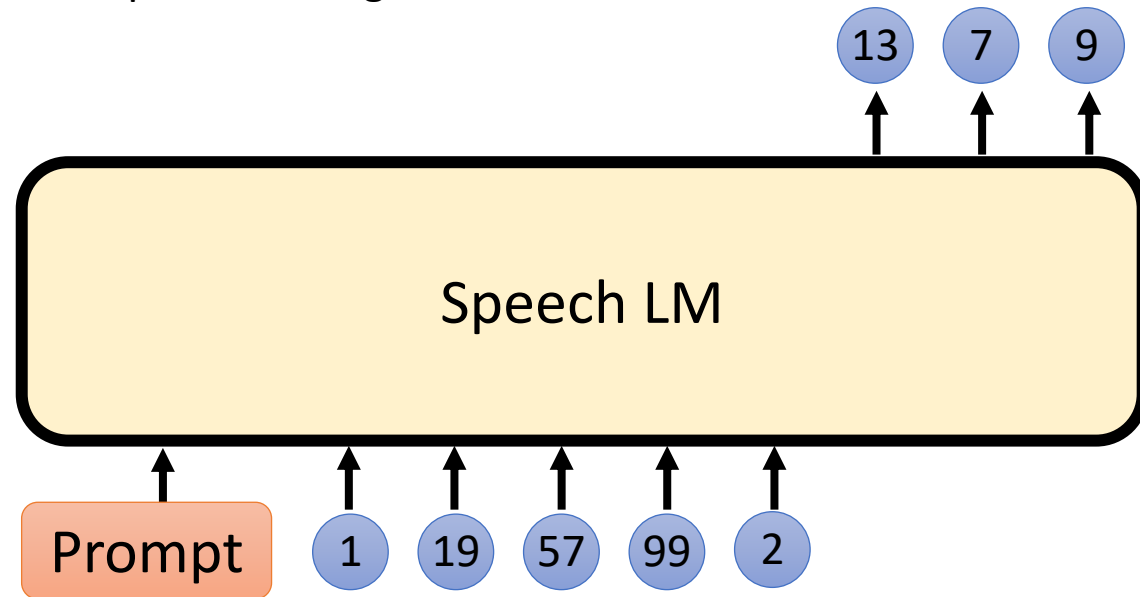


GSLM

<https://arxiv.org/abs/2102.01192>

Audio LM

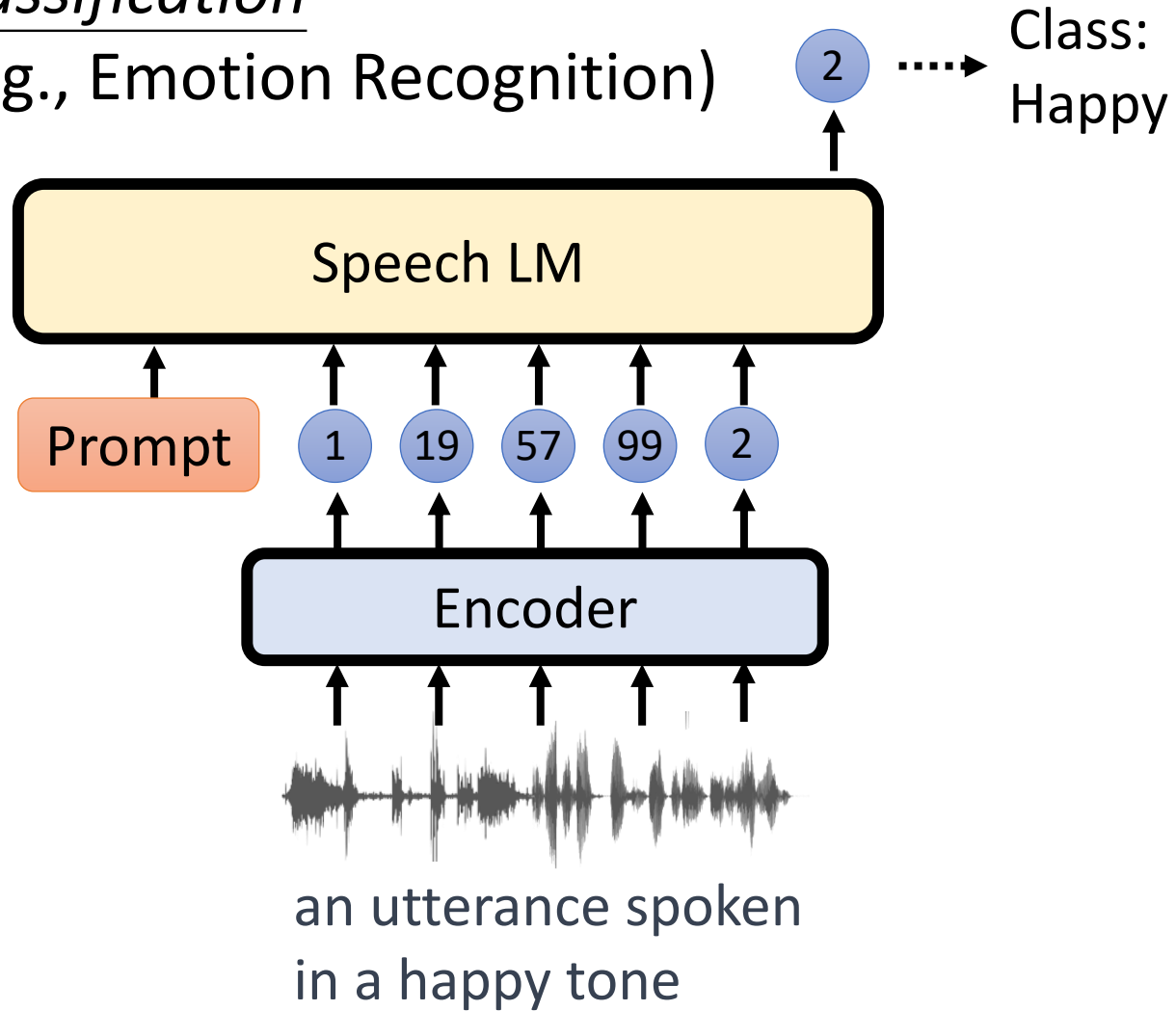
<https://arxiv.org/abs/2209.03143>



Does a speech LM have potential to address for a wide range of speech tasks?

Classification

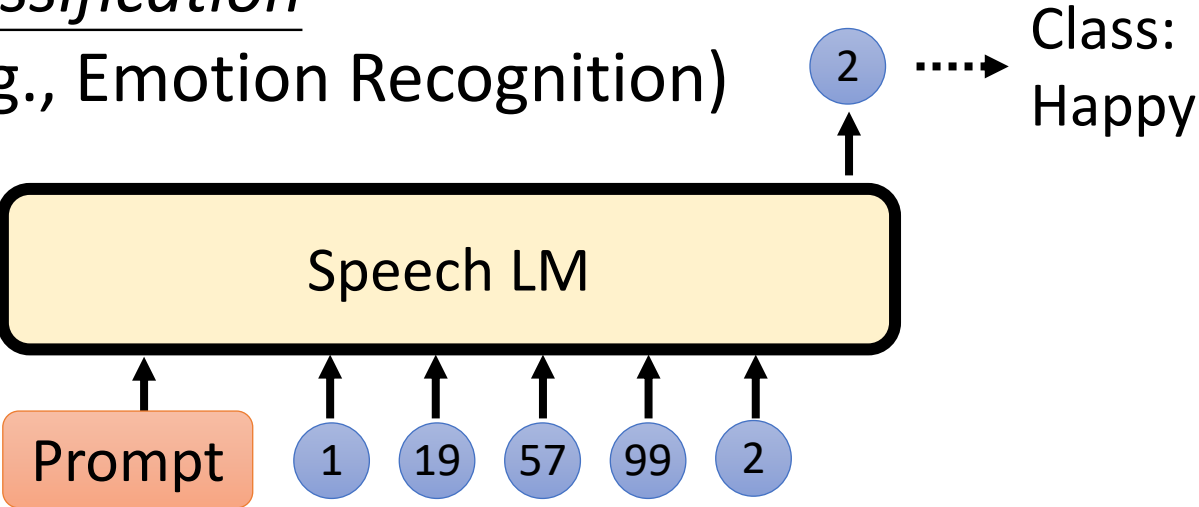
(e.g., Emotion Recognition)



Unit ID	Class
1	Angry
2	Happy
3	Sad
...	

Classification

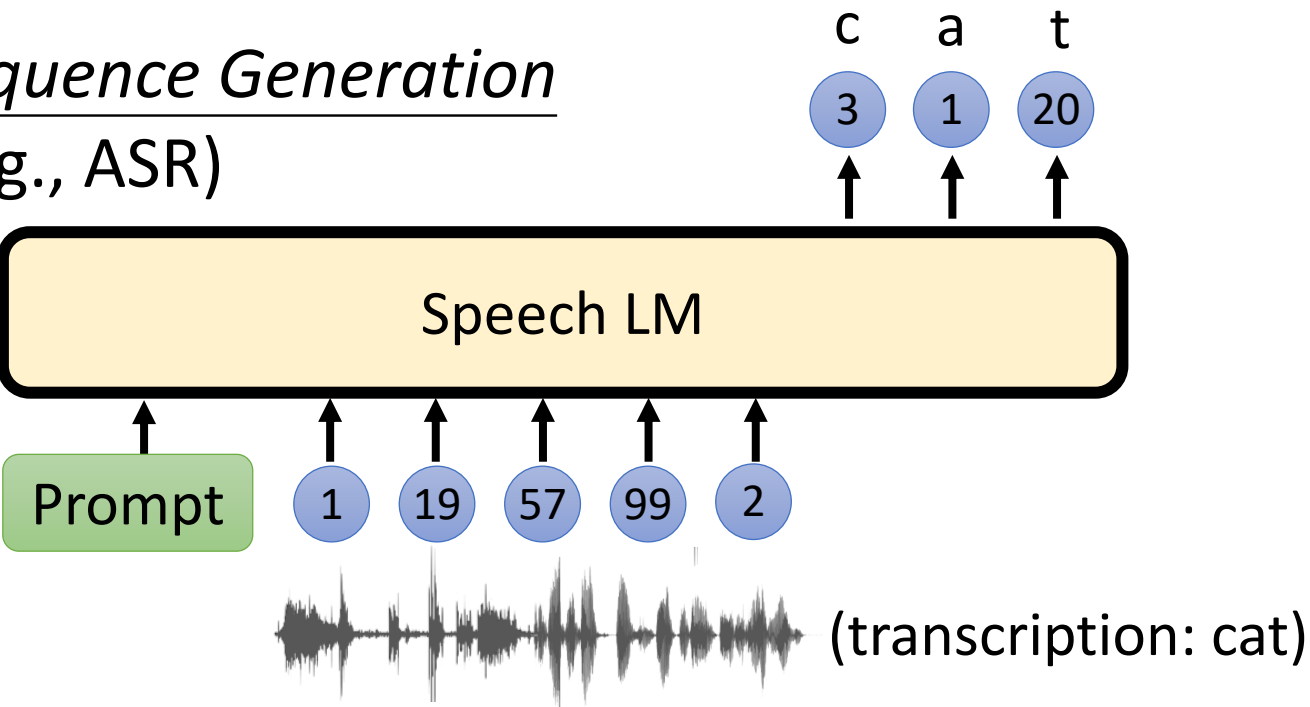
(e.g., Emotion Recognition)



Unit ID	Class
1	Angry
2	Happy
3	Sad
...	

Sequence Generation

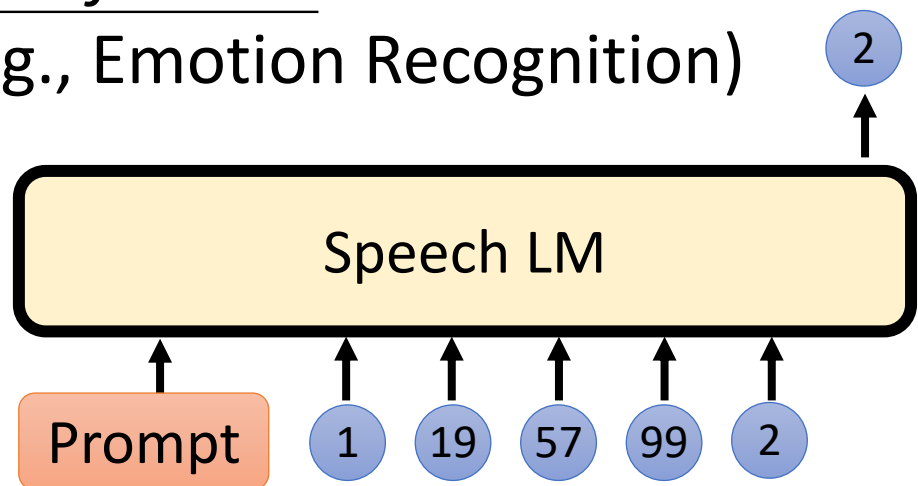
(e.g., ASR)



Unit ID	Class
1	a
2	b
3	c
...	

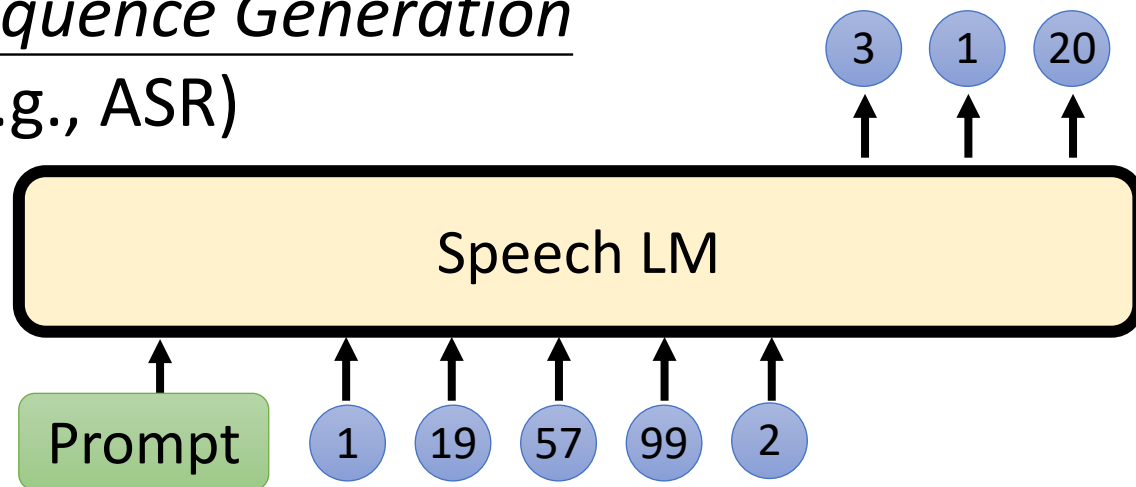
Classification

(e.g., Emotion Recognition)



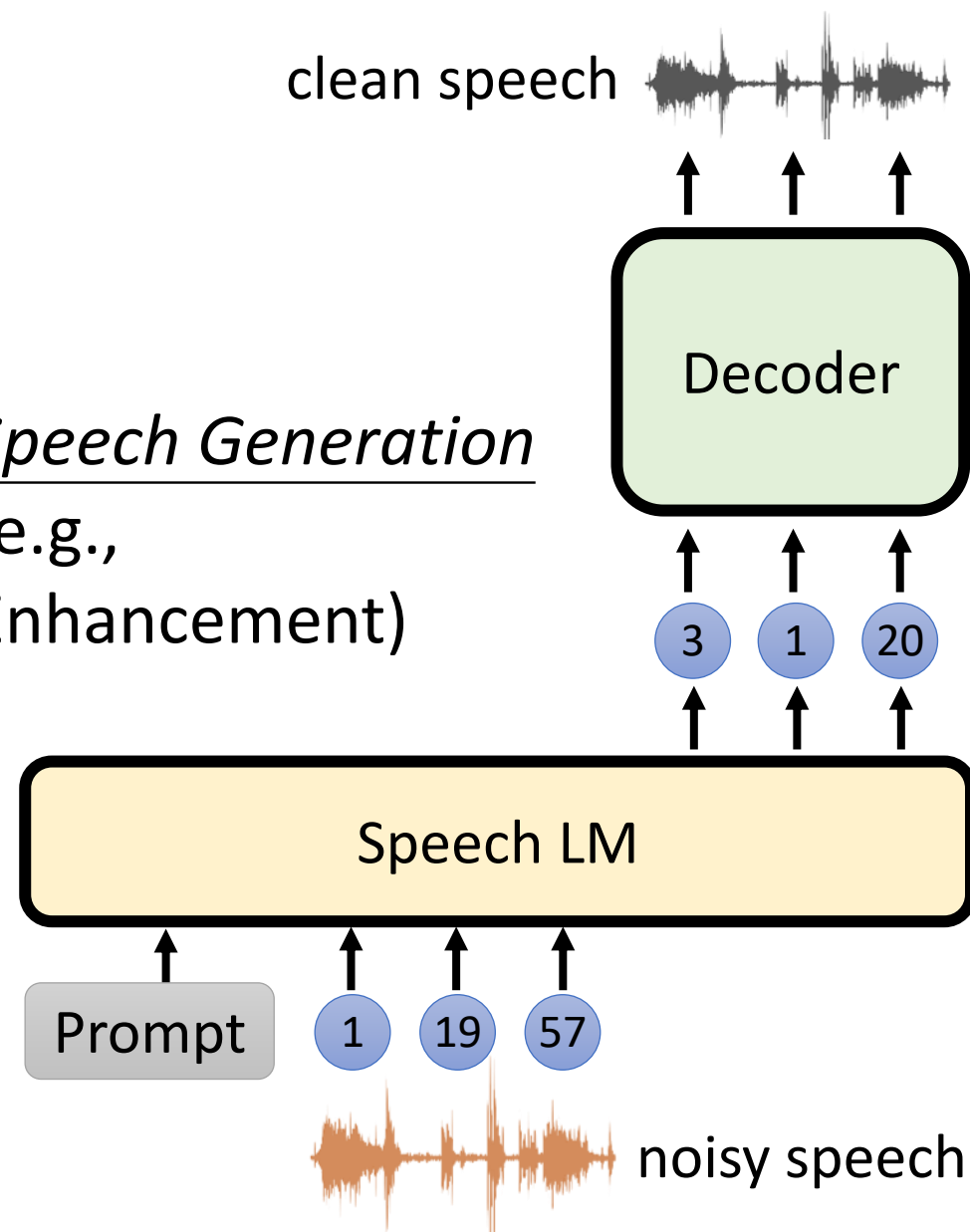
Sequence Generation

(e.g., ASR)



Speech Generation

(e.g., Enhancement)

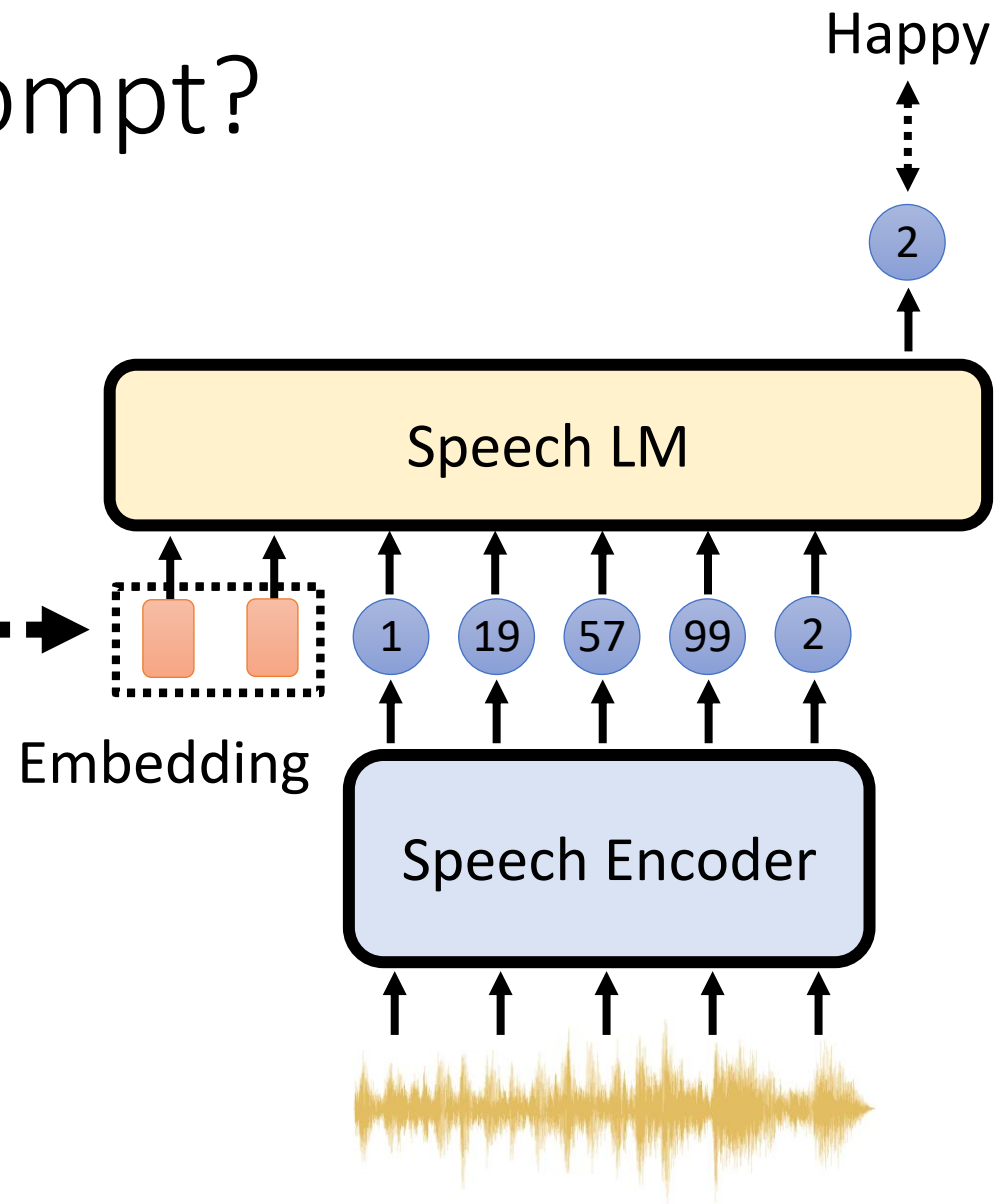


How to find the prompt?

Classification

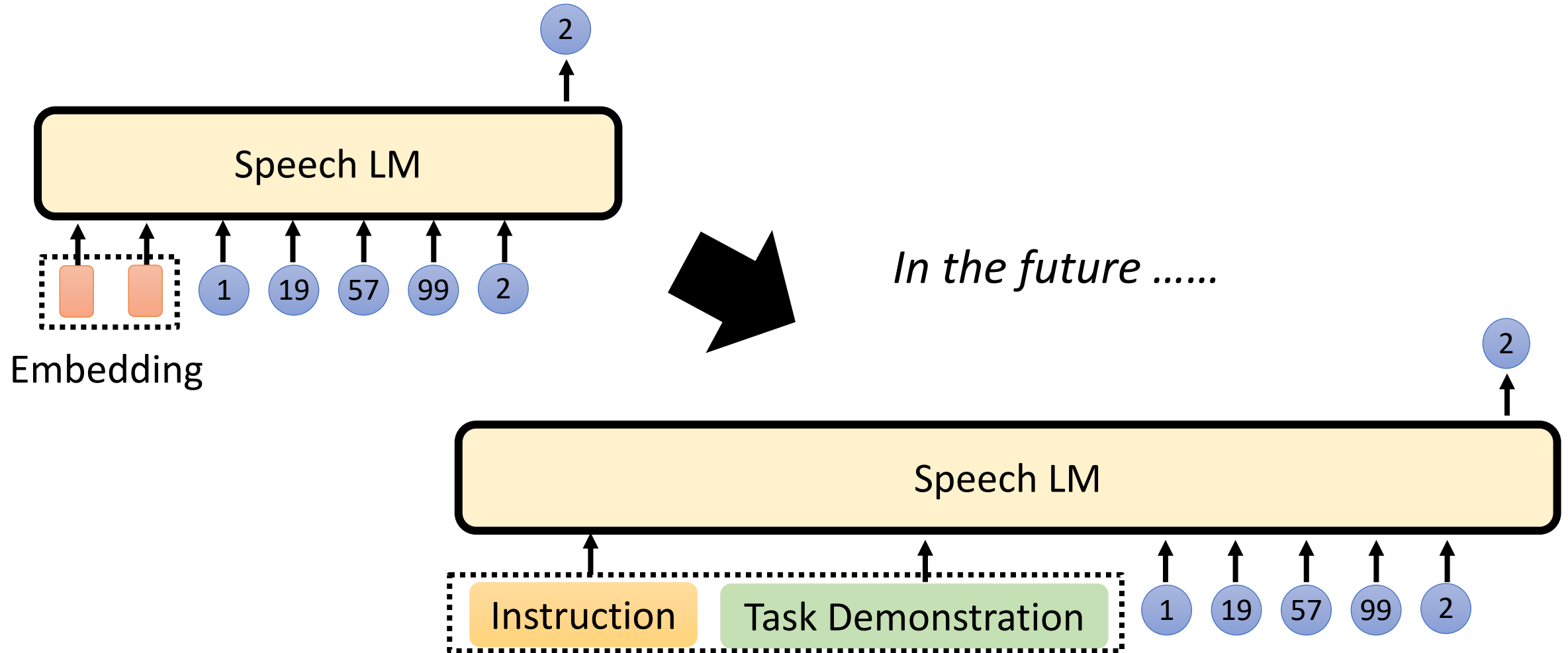
(e.g., Emotion Recognition)

we can learn prompt by
gradient descent



Unit ID	Class
1	Angry
2	Happy
3	Sad
...	

How far are we from Universal Speech Model?



SPEECHPROMPT V2: PROMPT TUNING FOR SPEECH CLASSIFICATION TASKS

Kai-Wei Chang^{1}*

Yu-Kai Wang^{2}*

Hua Shen³

Iu-thing Kang⁴

Wei-Cheng Tseng¹

Shang-Wen Li⁵

Hung-yi Lee¹

<https://arxiv.org/abs/2303.00733>



Kai-Wei Chang



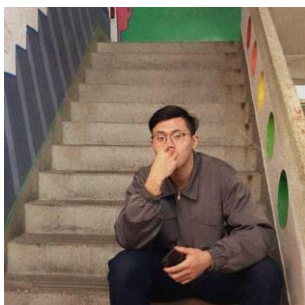
Yu-Kai Wang



Hua Shen



Iu-thing Kang



Wei-Cheng Tseng



Shang-Wen Li



Hung-yi Lee

Speech Classification Tasks

Speech Command
Recognition

Intent
Classification

Language
Identification

Dysarthric
Classification

Gender
Identification

Speaker
Identification

Accent
Classification

Sarcasm
Detection

Emotion
Recognition

Fake Speech
Detection

Voice Activity
Detection

Audio
Classification



GSLM
pGSLM



Prompt



+:Trainable mapping

Task	Metric	Dataset	Language	#Class	SOTA (Topline)	GSLM	GSLM+	pGSLM	pGSLM+
SCR	ACC (↑)	Google SC v1	En	12	98.6 [10]	94.5	94.6	94.3	94.7 (-3.9)
		Grabo SC	Du	36	98.9 [11]	92.4	92.7 (-6.2)	17.5	19.6
		Lithuanian SC	Lt	15	91.8 [9]	93.2	95.5 (+3.7)	90.9	79.5
		Arabic SC	Ar	16	98.9 [9]	99.7	100.0 (+1.1)	85.6	92.6
IC	ACC (↑)	Fluent SC	En	24	99.7 [12]	97.2	97.3	98.1	98.2 (-1.5)
LID	ACC (↑)	Voxforge	En, Es, Fr De, Ru, It	6	99.8 [13]	90.9	94.2 (-5.6)	81.8	80.4
FSD	EER (↓)	ASVspooF	En	2	2.5 [13]	18.5	13.5	13.1 (+10.6)	18.3
ER	ACC (↑)	IEMOCAP	En	4	79.2 [13]	42.1	44.3	49.9	50.2 (-29)
AcC	ACC (↑)	AccentDB	En	9	99.5 [14]	78.9	83.4	86.5	87.1 (-12.4)
SD	F1 (↑)	MUStARD	En	2	64.6 [15]	55.0	77.8	74.4	78.7 (+13.1)
		MUStARD++	En	2	65.2 [16]	74.0	75.2 (+10)	52.7	58.2
GID	F1 (↑)	VoxCeleb1	En	2	98.3 [17]	86.2	87.3	91.6 (-6.7)	86.2
VAD	ACC (↑)	Google SC v2 & Freesound	En	2	98.8 [18]	96.6	96.9	98.3 (-0.5)	98.1
AuC	ACC (↑)	ESC-50	✖	50	97.0 [19]	9.0	37.5 (-59.5)	20.3	27.0

We can get reasonable performance on most classification tasks with prompting.

Experimental Results – Sequence Generation

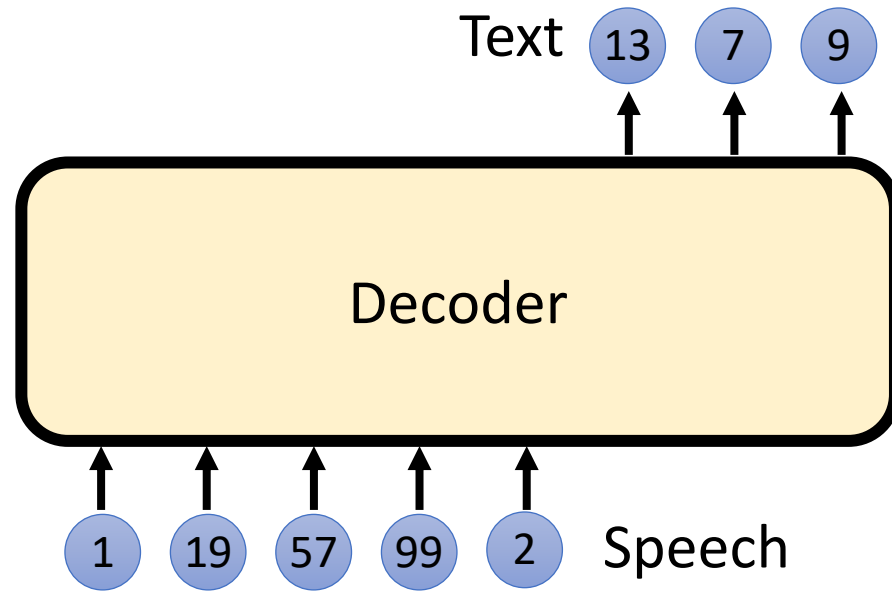
Prompt length = 180

		Slot Filling			
		ASR		SF	
Scenarios		WER ↓	CER ↓	F1 ↑	CER ↓
Prompt Speech LM	HuBERT-PT	34.17	26.14	66.90	59.47
Fine-tune Speech LM	FT-LM	26.19	16.80	80.58	40.15
SUPERB setting	FT-DM	6.42	1.48	88.53	25.20
Prompt Speech LM	CPC-PT	59.41	37.12	65.25	60.84
Fine-tune Speech LM	FT-LM	35.61	17.90	79.34	42.64
SUPERB setting	FT-DM	20.18	5.25	71.19	49.91

Prompting speech LM is worse than other approaches

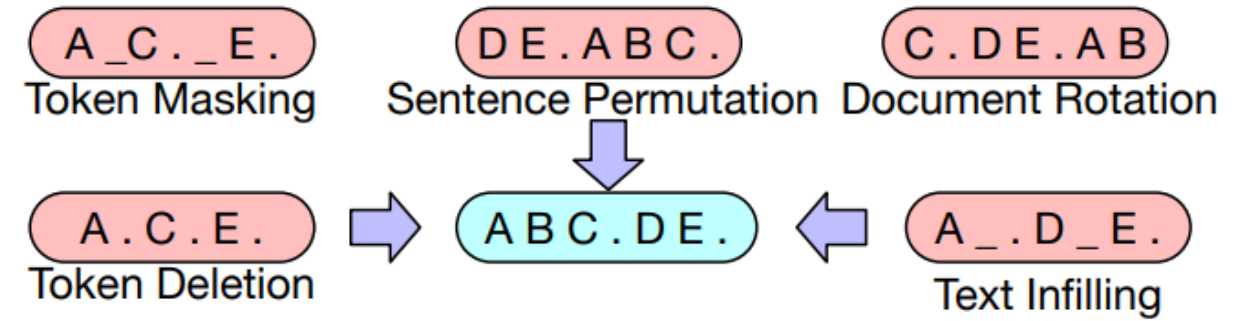
Other Speech LM

Poor performance for ASR

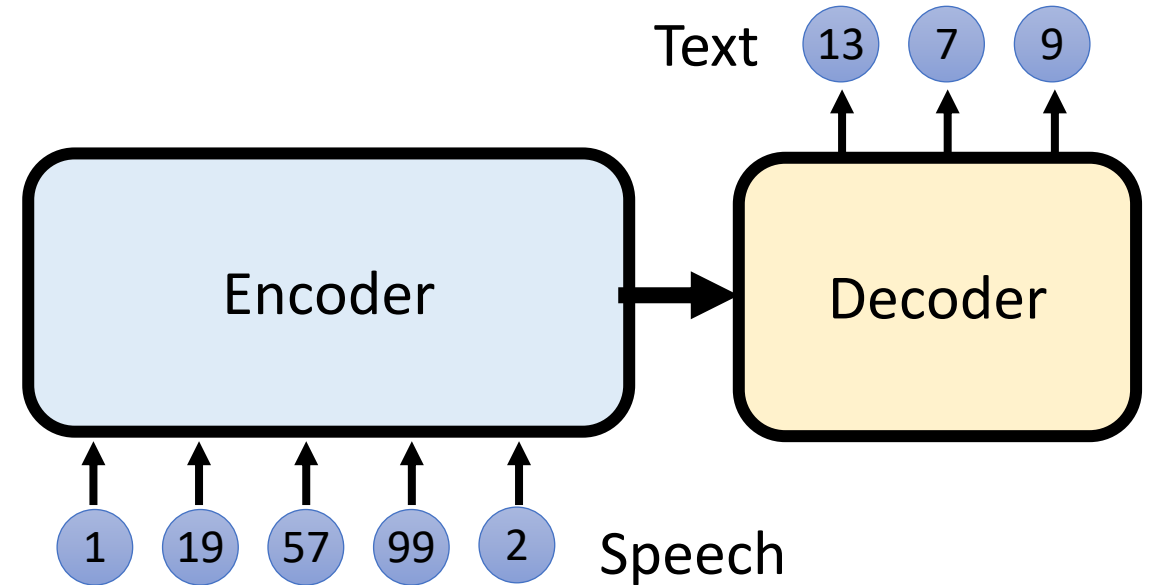


GSLM

<https://arxiv.org/abs/2102.01192>



<https://arxiv.org/abs/1910.13461>



Unit BART

<https://arxiv.org/abs/2204.02967>

Experimental Results – Sequence Generation

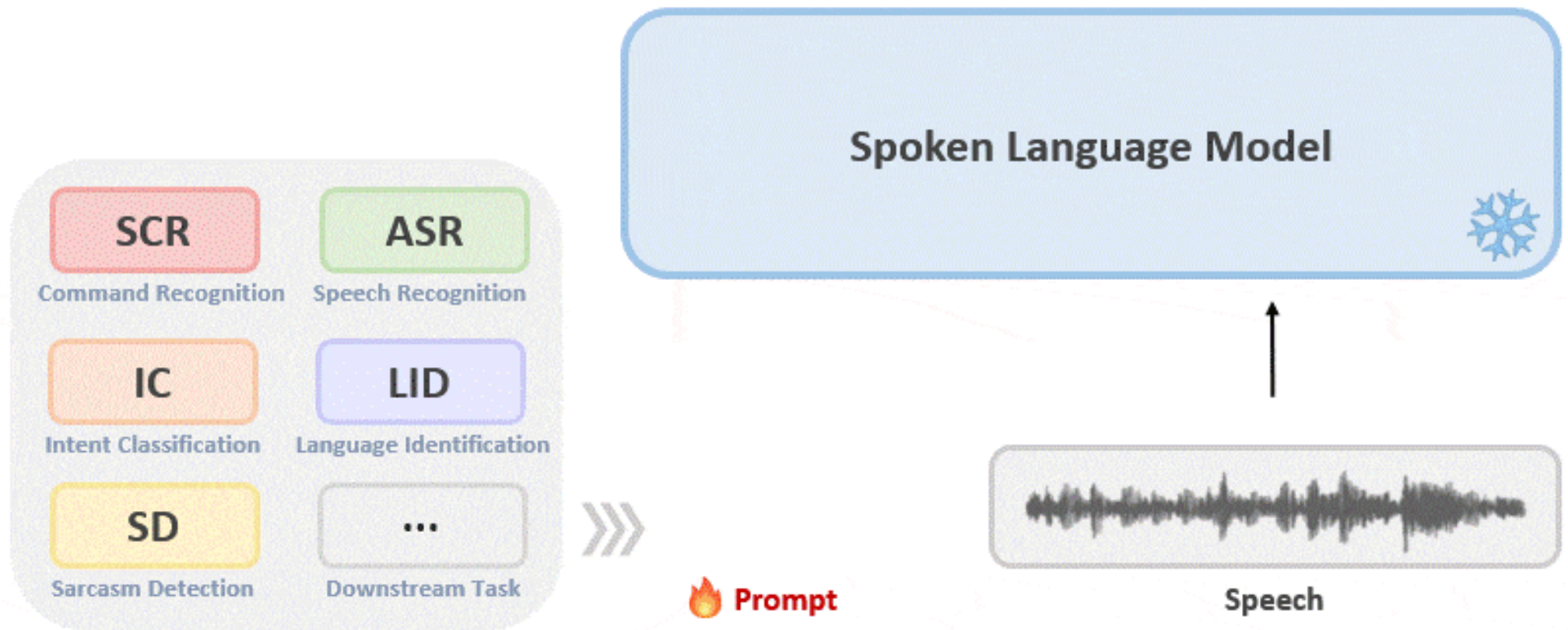
Slot Filling					
Scenarios	ASR		SF		#
	WER ↓	CER ↓	F1 ↑	CER ↓	
Prompt Speech LM→ HuBERT-PT	34.17	26.14	66.90	59.47	4.5M
Fine-tune Speech LM→ FT-LM	26.19	16.80	80.58	40.15	151M
SUPERB setting→ FT-DM	6.42	1.48	88.53	25.20	43M
Prompt Speech LM→ CPC-PT	59.41	37.12	65.25	60.84	4.5M
Fine-tune Speech LM→ FT-LM	35.61	17.90	79.34	42.64	151M
SUPERB setting→ FT-DM	20.18	5.25	71.19	49.91	42.5M

Prompting unit BART obtains 9.8% CER.

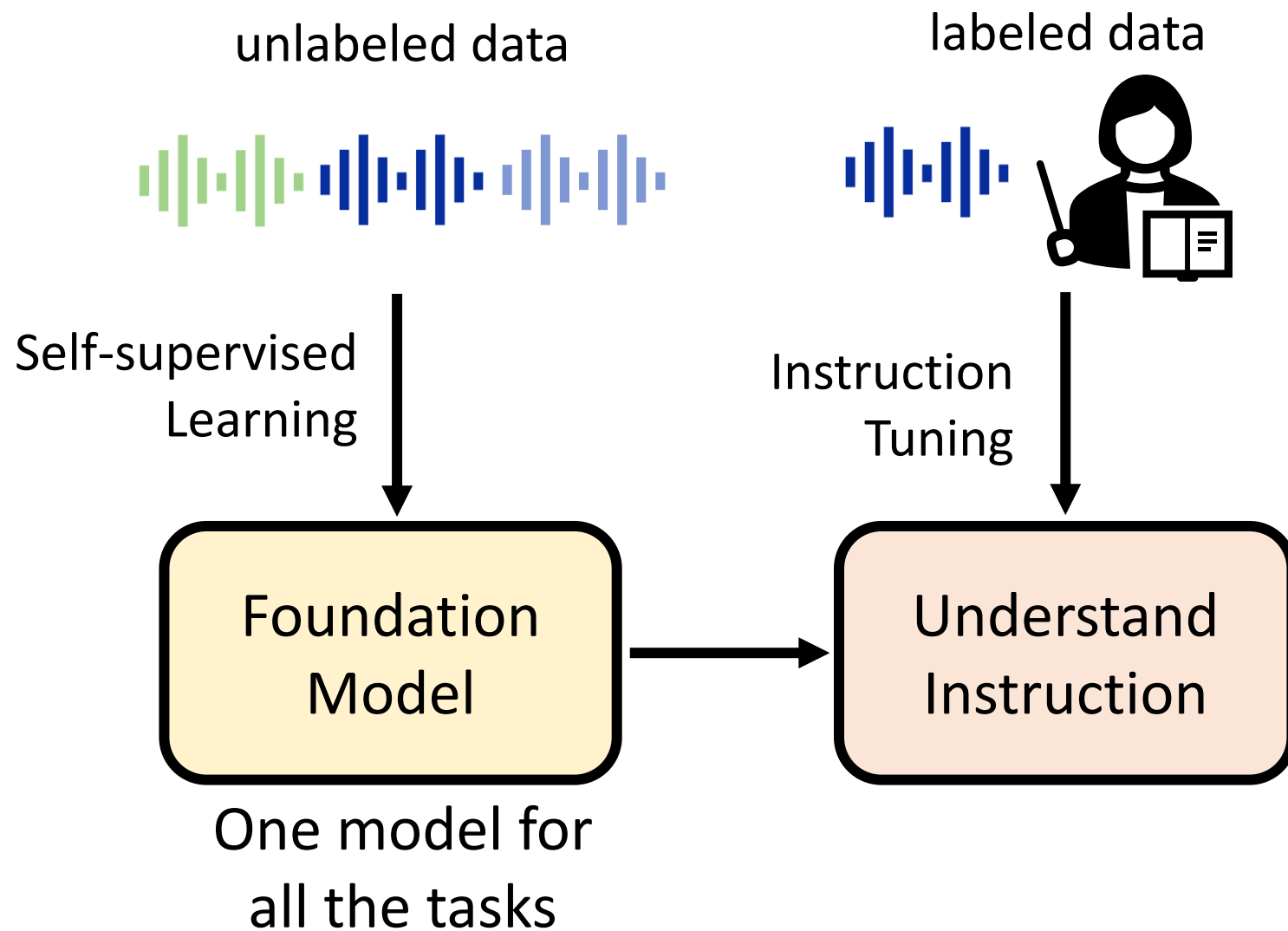
(unpublished results)

To learn more

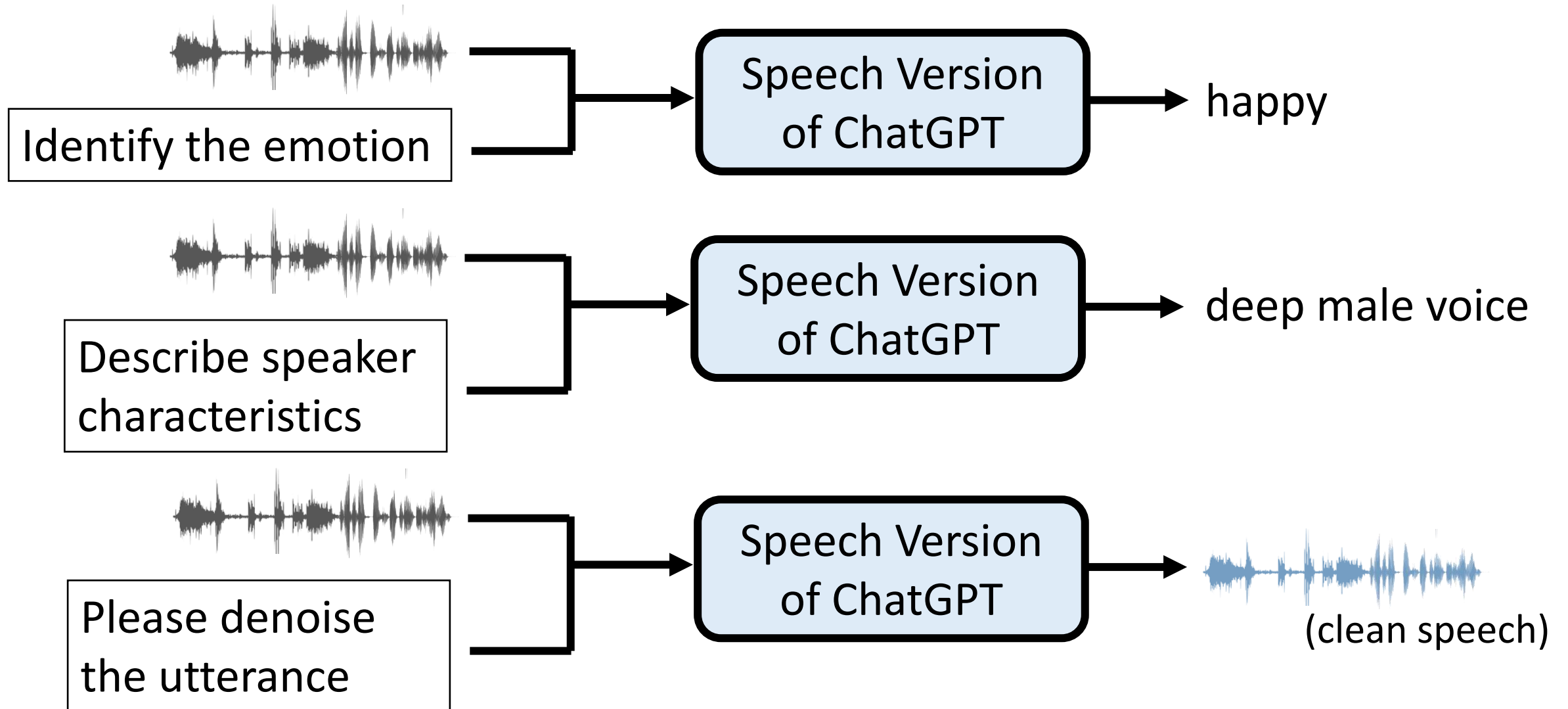
<https://ga642381.github.io/SpeechPrompt/>



How about Speech?



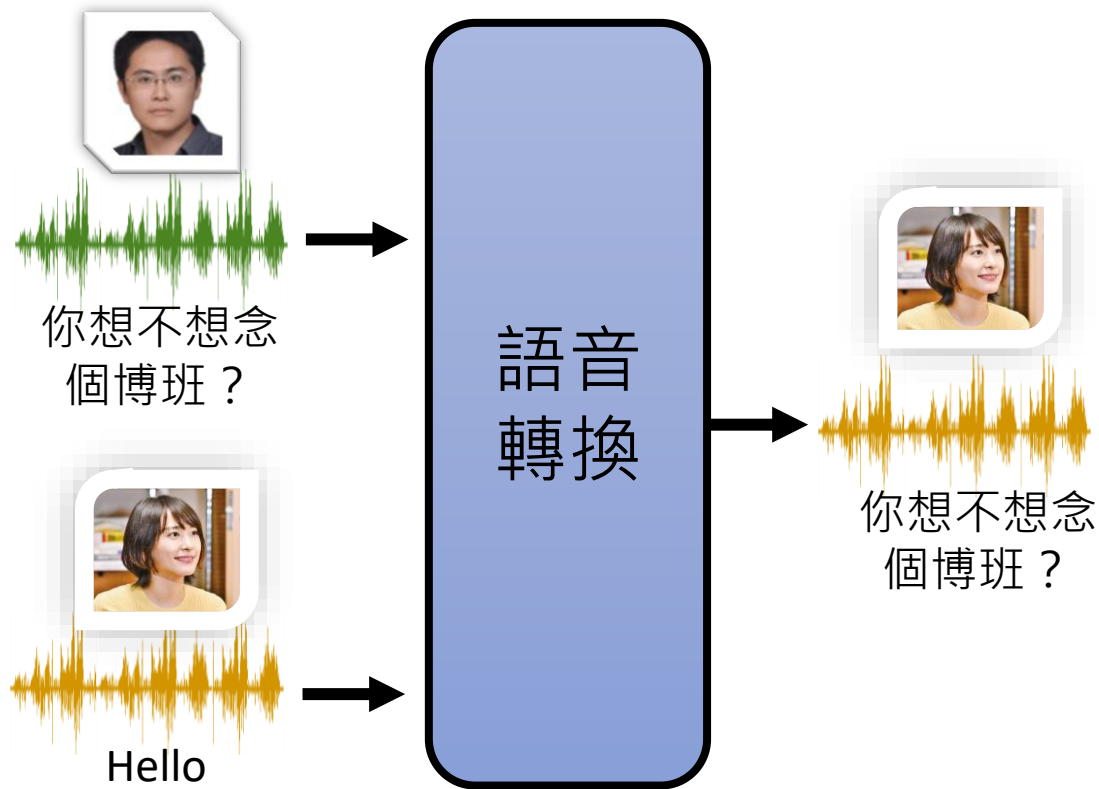
Speech Version of ChatGPT



Instruction Fine-tuning for Speech LM: Text-instruction-guided Voice Conversion



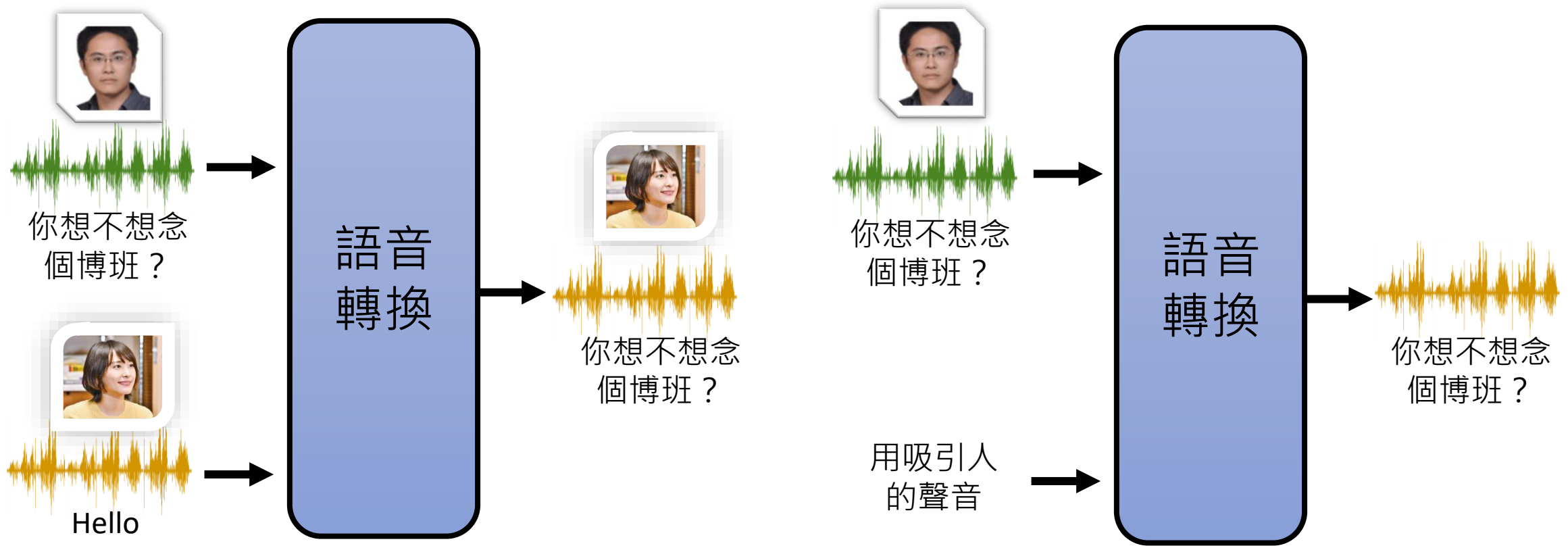
Instruction Fine-tuning for Speech LM: Text-instruction-guided Voice Conversion



INTERSPEECH 2022 Tutorial
(Xu Tan and Hung-yi Lee)

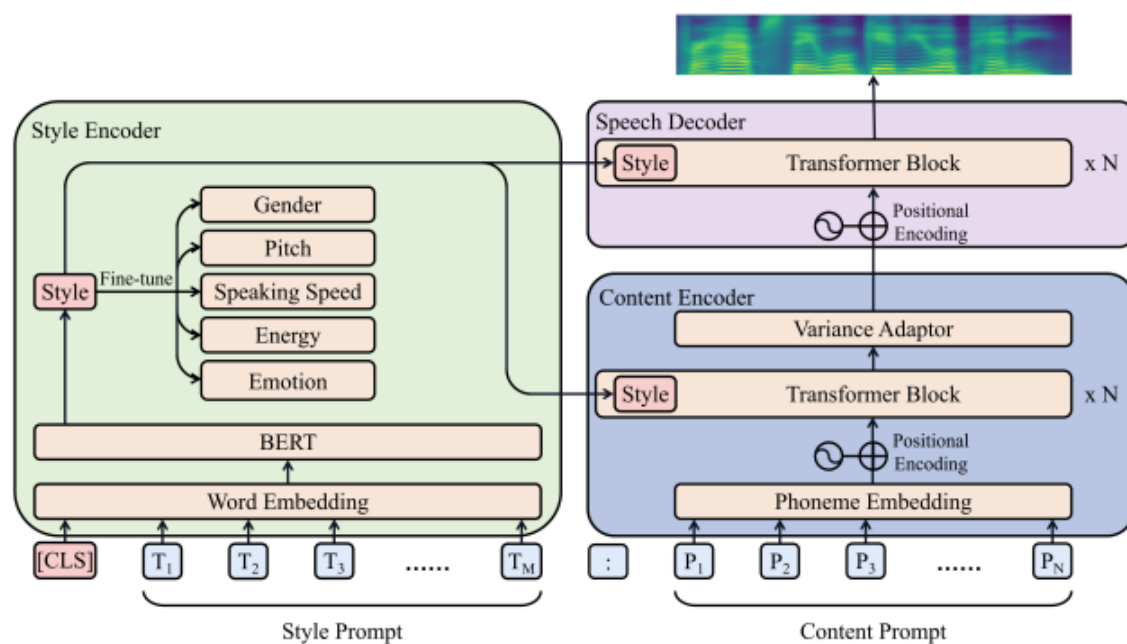
<https://tts-tutorial.github.io/interspeech2022/>

Instruction Fine-tuning for Speech LM: Text-instruction-guided Voice Conversion



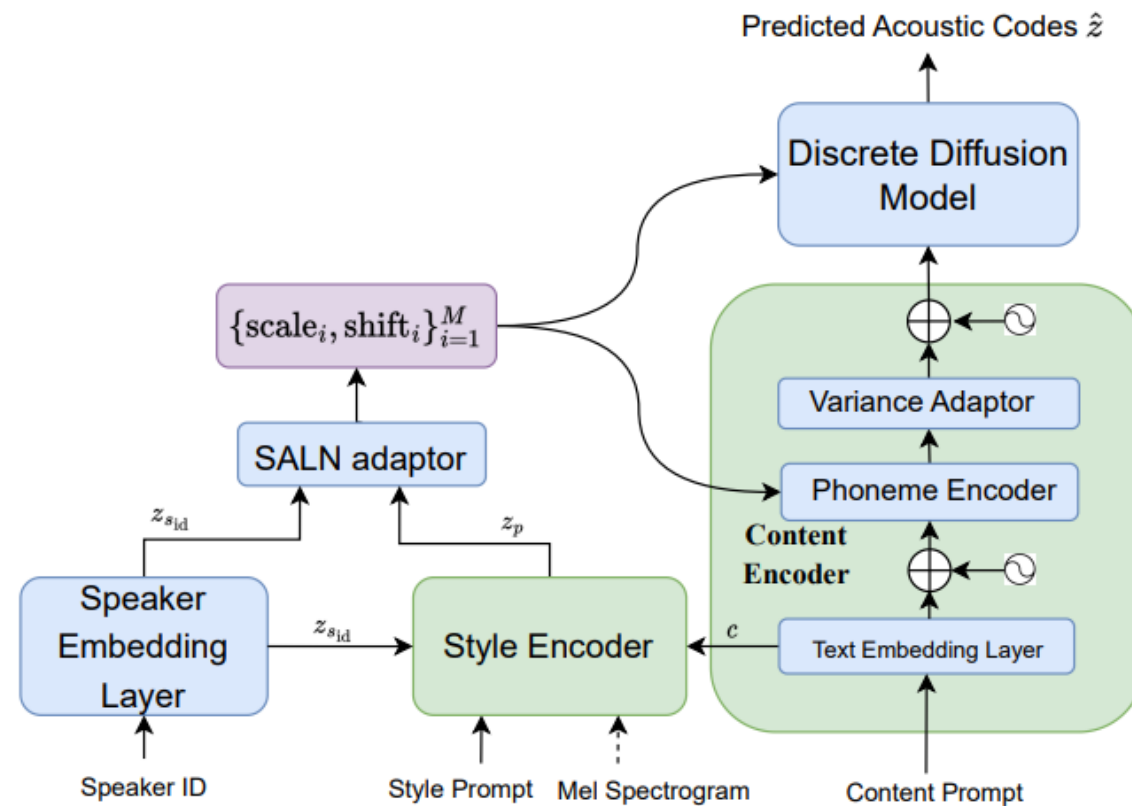
<https://youtu.be/JZvEzb5PV3U>

Related Work



PromptTTS

<https://arxiv.org/abs/2211.12171>



InstructTTS

<https://arxiv.org/abs/2301.13662>

Instruction Fine-tuning for Speech LM: Text-instruction-guided Voice Conversion



Chun-Yi Kuan
(NTU)



Chen-An Li
(NTU)



Tsu-Yuan Hsu
(NTU)



Tse-Yang Lin
(NTU)



Ho-Lam Chung
(NTU)



Kai-Wei Chang
(NTU)

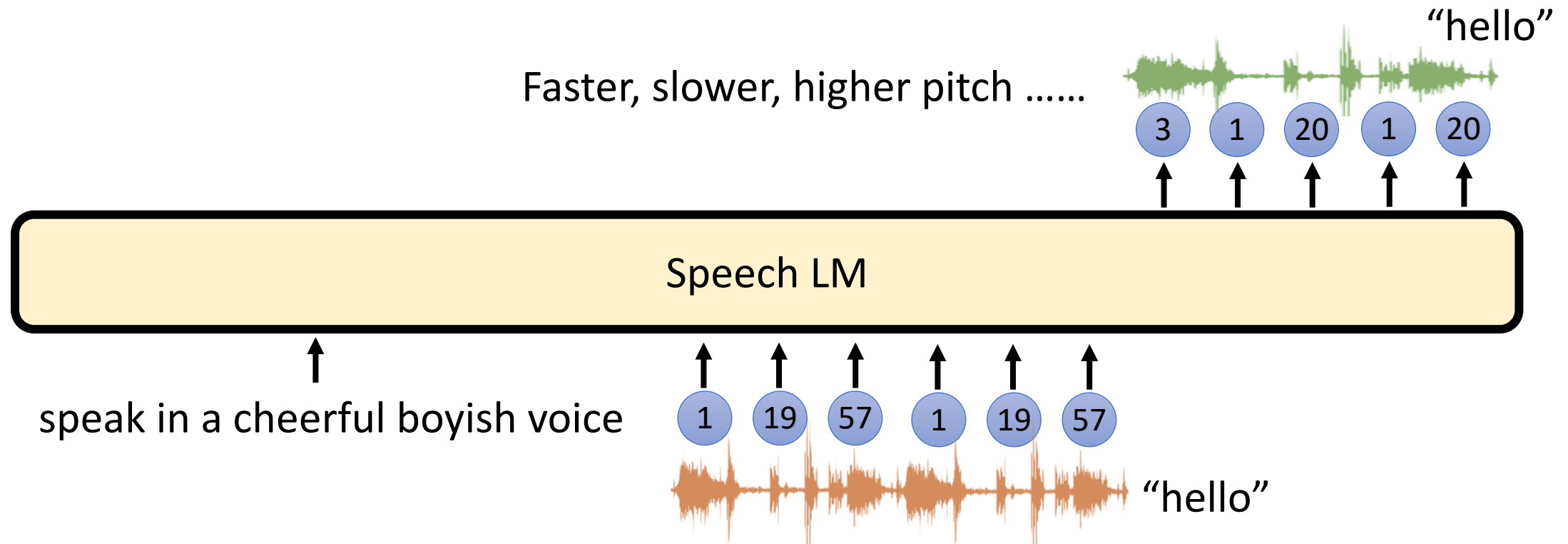


Shuo-yiin Chang
(Google)

This research is based on work that is funded by a Googler-Initiated Grant from Shuo-Yiin Chang.

Instruction Fine-tuning for Speech LM

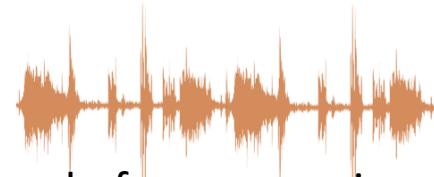
- Preliminary study: Text-instruction-guided Voice Conversion



Text-instruction-guided Voice Conversion

Training data format:

Text Instruction



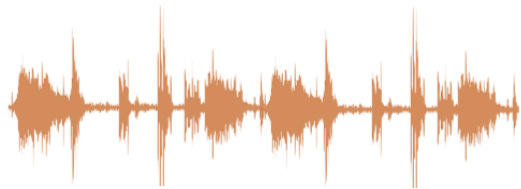
before conversion



after conversion

Signal Processing Effect Dataset

sox -G src.wav tgt.wav **tempo** 1.75



Template + ChatGPT

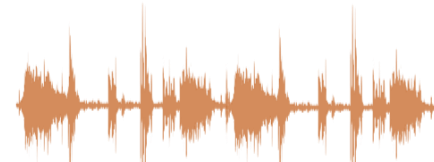
Speed up the audio to a significant degree.



Text-instruction-guided Voice Conversion

Training data format:

Text Instruction

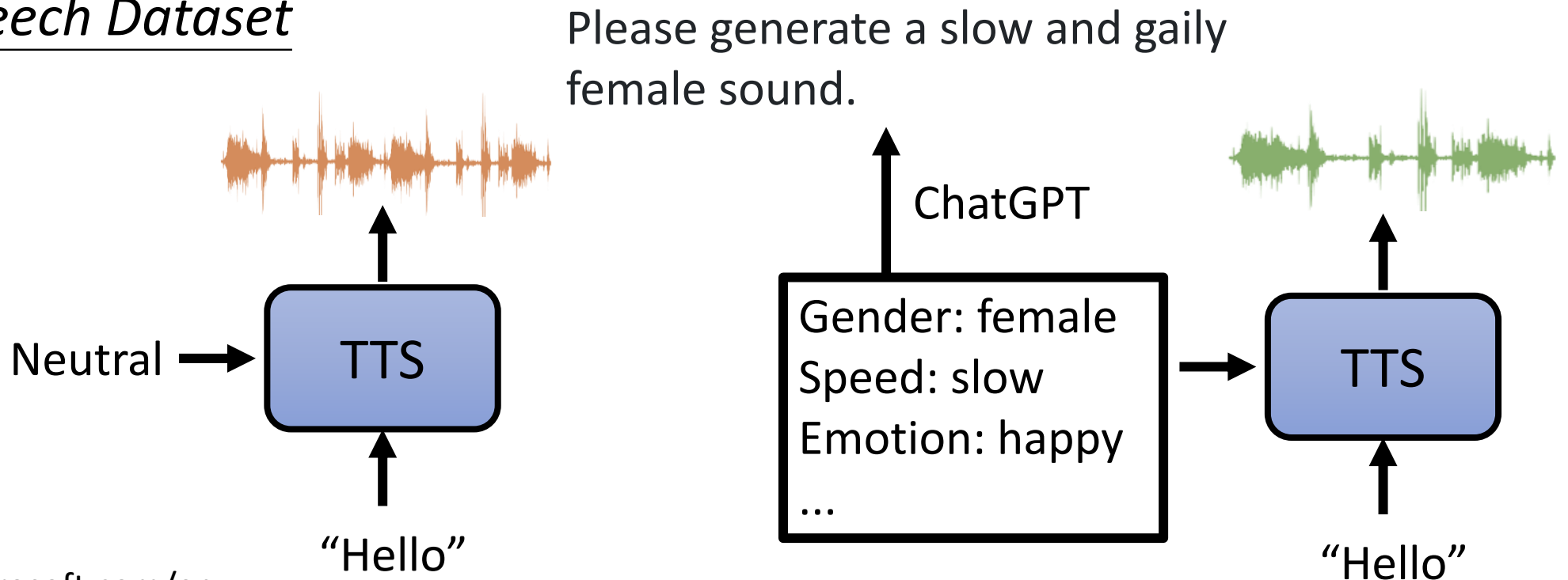


before conversion



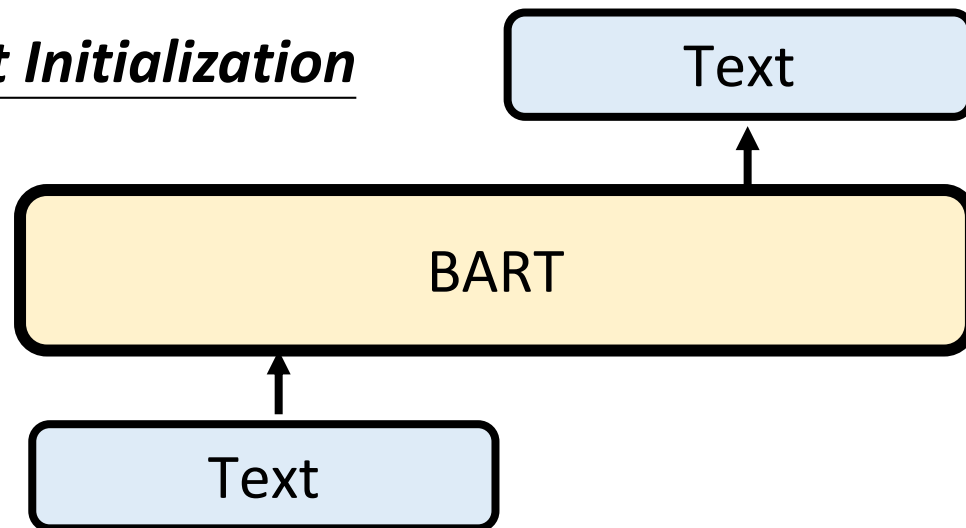
after conversion

InstructSpeech Dataset

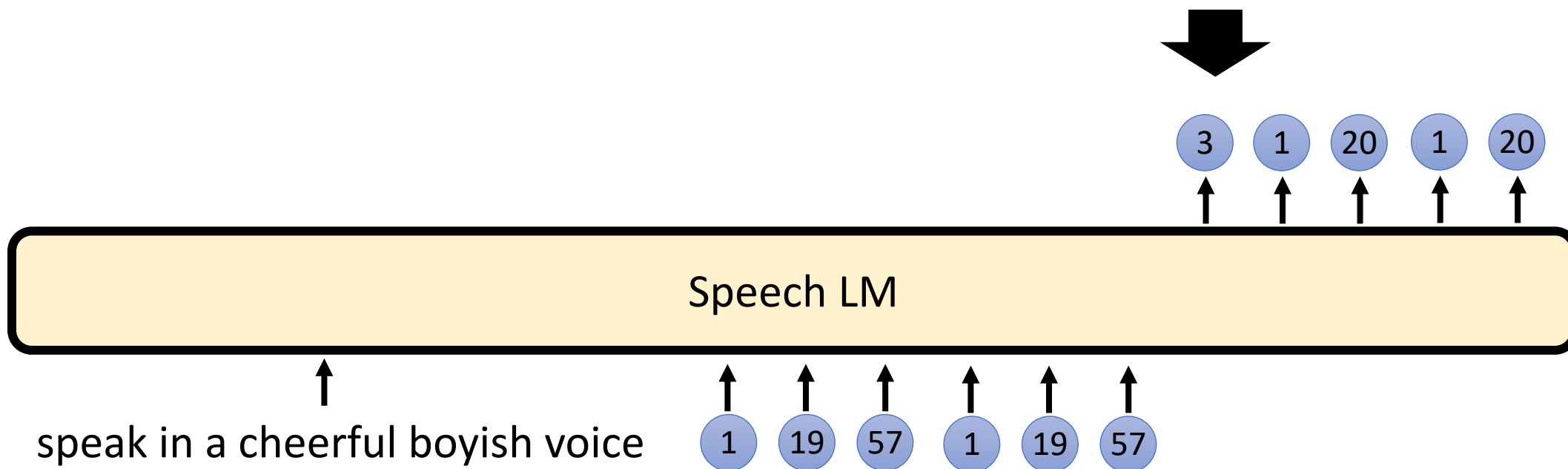
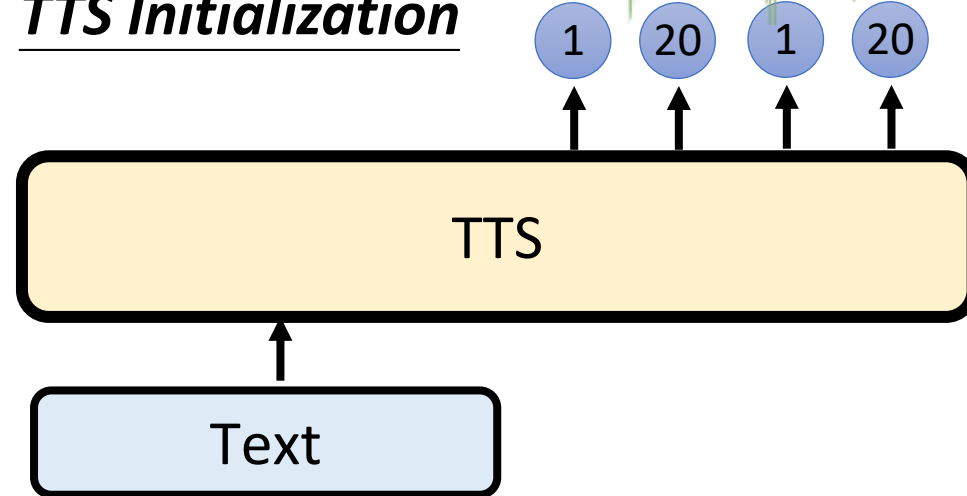


TTS: <https://azure.microsoft.com/en-us/products/ai-services/text-to-speech>

Text Initialization



TTS Initialization



Experimental Results

Human Evaluation (MOS Score)

	Quality	Instruction
Text-guided VC	2.81	4.19
Ground Truth	4.15	3.69

- Don't worry. We have improved the quality by using a better EnCodec decoder.
- Even better than ground truth? The model generates more apparent conversion.

Experimental Results

A slow sad moving female bass appeared in a low volume.



Add a profound sense of spatial dimension for a more immersive audio experience.



Speak as if you're telling a bedtime story to a child. **(no similar training instructions)**



Adopt the tone of a news anchor delivering breaking news. **(no similar training instructions)**



Text Pre-
training

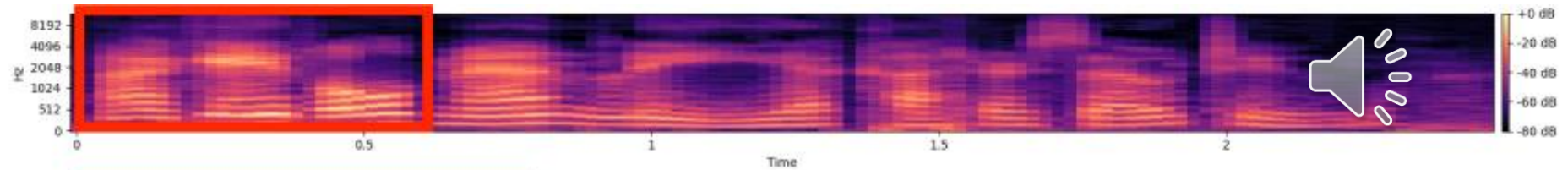


No Pre-
training

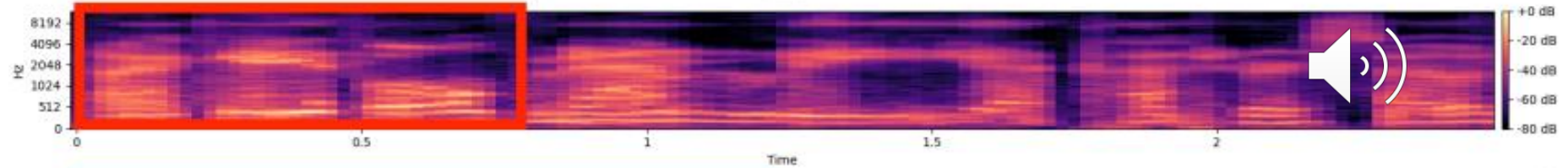


Experimental Results

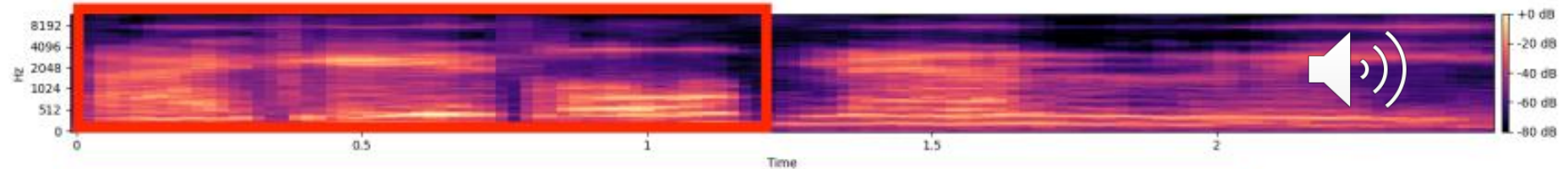
Source speech



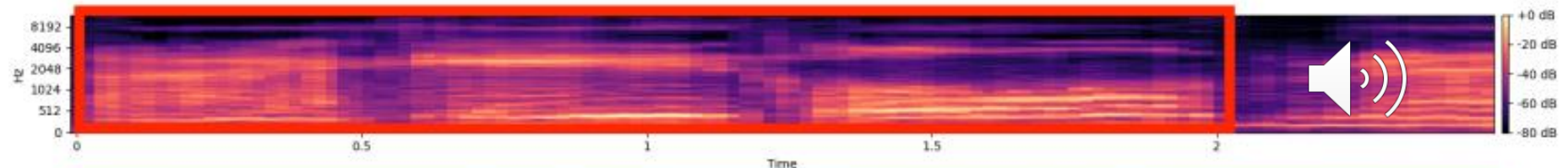
Decrease the speed of speech **slightly**.



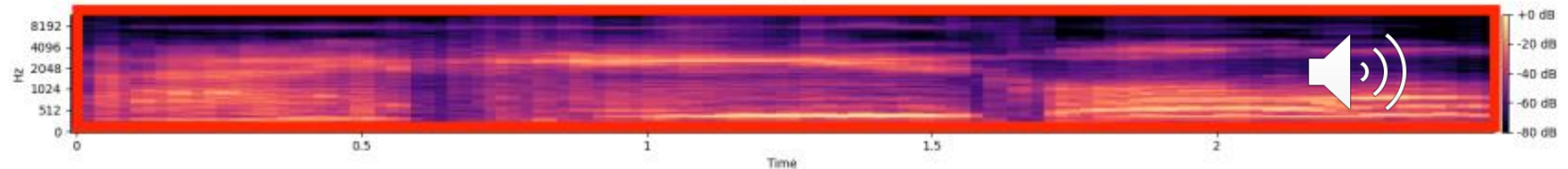
Decrease the speed of speech.



Decrease the speed of speech **notably**.



Decrease the speed of speech **extremely**.



What is still missing?

Speech

SUPERB series: 17 tasks

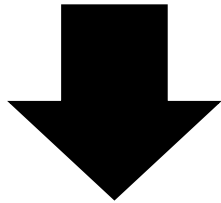
<https://arxiv.org/abs/2105.01051>

<https://arxiv.org/abs/2203.06849>

<https://arxiv.org/abs/2210.08634>

<https://arxiv.org/abs/2210.07185>

<https://arxiv.org/abs/2110.06280>



We need something bigger.

with instruction 😊

NLP

General Language Understanding
Evaluation (GLUE): 9 tasks

<https://arxiv.org/abs/1804.07461>

Super GLUE: 8 tasks

<https://arxiv.org/abs/1905.00537>

FLAN: 62 tasks

<https://arxiv.org/abs/2109.01652>

CrossFit: 160 tasks

<https://arxiv.org/abs/2104.08835>

BIG-bench: 204 tasks

<https://arxiv.org/abs/2206.04615>

natural-instructions: 1616 tasks

<https://arxiv.org/abs/2204.07705>

DYNAMIC-SUPERB: TOWARDS A DYNAMIC, COLLABORATIVE, AND COMPREHENSIVE INSTRUCTION-TUNING BENCHMARK FOR SPEECH

*Chien-yu Huang¹, Ke-Han Lu^{*1}, Shih-Heng Wang^{*1}, Chi-Yuan Hsiao^{†1}, Chun-Yi Kuan^{†1}, Haibin Wu^{†1}
Siddhant Arora^{§2}, Kai-Wei Chang^{§1}, Jiatong Shi², Yifan Peng², Roshan Sharma², Shinji Watanabe²
Bhiksha Ramakrishnan^{2,3}, Shady Shehata³, Hung-yi Lee¹*

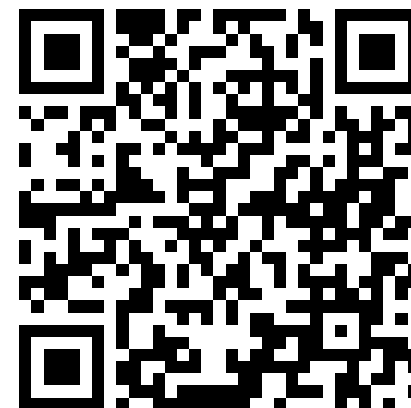
¹National Taiwan University, Taiwan, ²Carnegie Mellon University, USA

³Mohamed bin Zayed University of Artificial Intelligence, United Arab Emirates

<https://arxiv.org/abs/2309.09510>





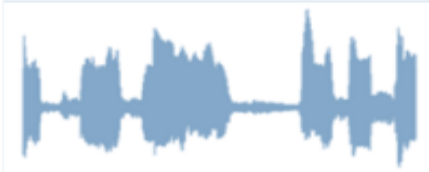
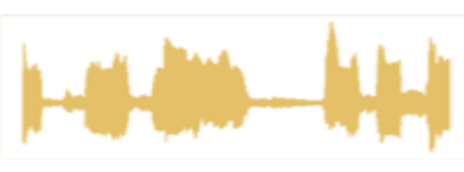


Chien-yu
Huang



Project page: <https://github.com/dynamic-superb/dynamic-superb>

Format

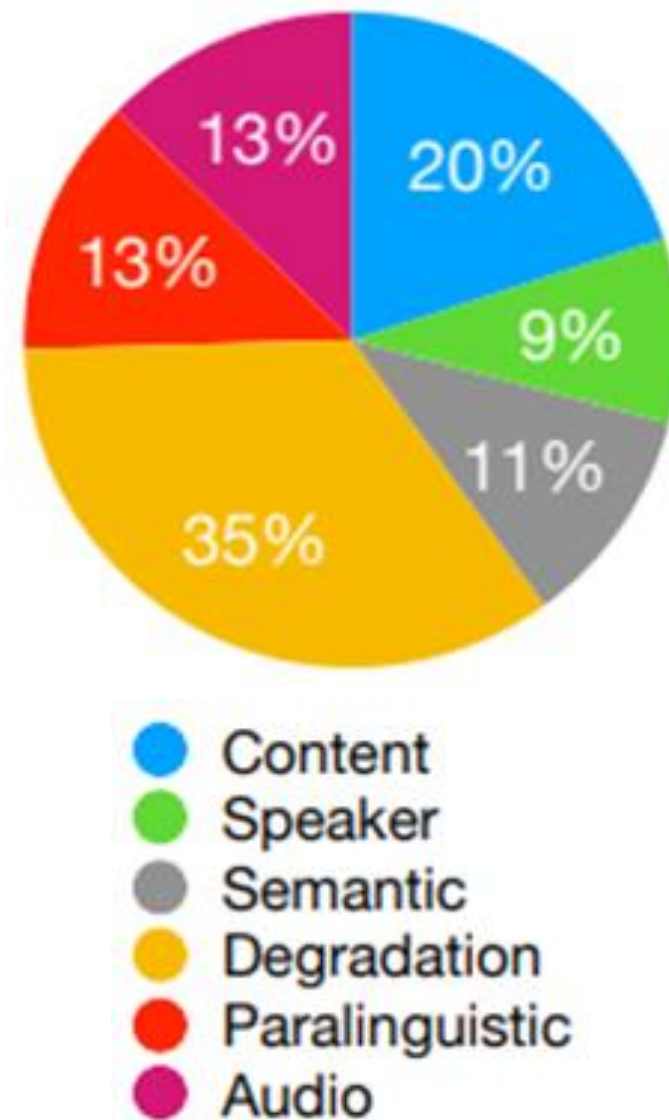
Instruction	Input	Output
Please identify the emotion in the audio. The answer could be		“Happiness”
Assess whether the provided speech is a spoofed voice. The answer could be		“Spoofed”
Recognize when sarcasm or irony is employed in the speech. The answer could be		“True”
Please transcribe the utterance.		“How are you”
Please speak slower.		

Current Status

- 55 tasks created from 33 datasets
- Covering 6 dimensions
 - Content: speech command recognition
 - Speaker: speaker verification
 - Semantics: sarcasm detection
 - Degradation: noise SNR prediction
 - Paralinguistic: emotion recognition
 - Audio: environmental sound classification

55 is not a big number ...

They are all classification tasks ...



Let's work together!

[illegible]

444 authors across 132 institutions

BEYOND THE IMITATION GAME: QUANTIFYING AND EXTRAPOLATING THE CAPABILITIES OF LANGUAGE MODELS

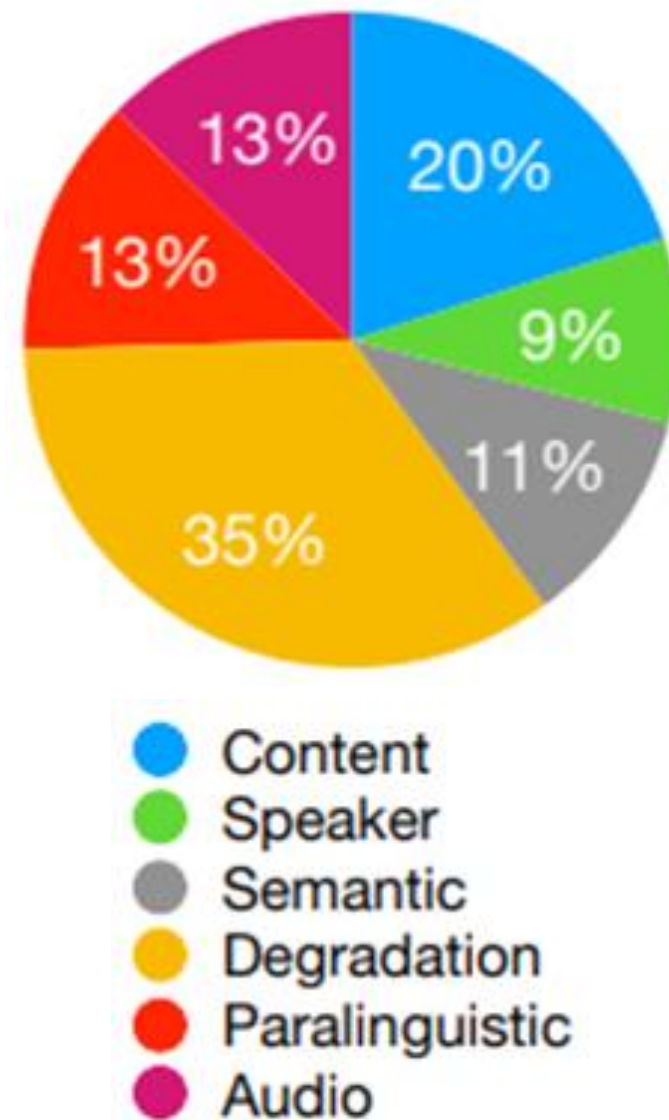
Aashu Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Asad Miri Shabbir, Abubakar Abdi, Adam Fisch, Adam R. Brown, Adam Samson, Aditya Gupta, Adria Garriga-Alonso, Agnieszka Kladka, Aitor Lemkowicz, Akshar Agarwal, Alina Pover, Alan Ray, Alex Watanabe, Alexander W. Kucuk, Ali Salata, Ali Tazari, Alicia Xiang, Alicia Parrish, Allen Nie, Amos Hassan, Amanda Asch, Amanda Dsouza, Ambrose Stone, Amotz Eshau, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Smalchiff, Andrew Dai, Andrew L. Andrew Lampinen, Andy Sun, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Goratchi, Antonia Nonell, Anna Vukobrat, Arash Ghahramani, Arfa Tabasum, Arif Muneer, Arun Kishan, Asher Malkoch, Ashish Sabharwal, Austin Herrick, Avia Eilat, Aykut Erdem, Ayta Karaday, B. Ryan Roberts, Bao Sheng Lee, Barret Zoph, Bartholomaj Bojanowski, Benjamin Ouyang, Behnam Heydari, Behnam Neyshabadi, Benjamin Peters, Benno Stein, Berk Elmekci, Bill Yuchen Lin, Blake Howald, Cameron Diao, Cameron Dour, Catherine Strauss, Cedrick Auguste, César Ferri Ramon, Chandan Singh, Charles Rathkopf, Chelsea Wang, Chitra Baral, Chiya Wu, Chir Caliskan-Burch, Chris Walter, Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Sirio, Colin Raffel, Cosma Schaefer, Cristina Garbacia, Damien Silve, Dan Garmon, Dan Hendrycks, Dan Kilian, Dan Roth, Daniel Freeman, Daniel Khazanchi, Daniel Levy, Daniel Mesquita Gonzalez, Danielle Portyck, Danny Hernandez, Dany Chen, Daphne Ippolito, Dar Gilboa, David Dolan, David Duval, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Denis Yuret, Derek Chen, Derek Tan, Desiree Hapkin, Dignata Miera, Dilyar Baran, Dimitri Costin-Mollo, Diyi Yang, Dong-Ho Lee, Ekaterina Shumova, Elia Dognas Gabali, Elad Segal, Elanor Hagmann, Elizabeth Barnes, Elizabeth Dunning, Elie Perle, Emanuele Rodola, Emma Lau, Eric Chu, Eric Tang, Erik Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Frazak, Ethan Kim, Ethan Fagola Manyasi, Evgeni Zhelenskiy, Fanyue Xia, Fatenah Sir, Fernando Martin-Pham, Francesca Happel, Francisco Chollet, Frieda Ring, Gaurav Mishra, Gema Duato Wiman, Gerard de Melo, Gerardo Kruczycki, Gianbattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovich-Lopez, Gregor Betz, Guy Gur-Ari, Hana Galjanovic, Hannah Kim, Hannah Raskin, Hannaneh Hajishirri, Harsh Mehta, Hayden Bogar, Henry Shevlin, Herich Schmitz, Hironori Yukawa, Hongming Zhang, Hugh Mao Wong, Ian Mi, Isaac Noble, Jaap-Janet, Jack Geisinger, Jackson Kervin, Jacob Hilton, Jacobson Lee, Jaime Fernandez Fdez, James B. Simon, James Koppel, James Zhang, James Zou, Jan Kozul, Jan Thompson, Janel Kanari, Janina Radom, Jascha Sobel-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bouchee, Jennifer Marsh, Jeremy Kim, Jesse Taal, Jesse Engel, Joseph Abadi, Kachun Xu, Jianing Song, Jifan Tang, Joan Waweru, John Barden, John Miller, John U. Balis, Jonathan Horvath, Jory Finkberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chan, Kamal Kancher, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Keja Markov, Kenneth D. Eboli, Kevin Gimpel, Kevin Omund, Kory Mathewson, Krzysztof Chatallo, Ksenia Shkara, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Lara Reynolds, Leo Gao, Li Zhang, Lian Dognas, Lianhui Qiu, Lina Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schleich, Ludwig He, Luis Oliveira Cohen, Luis Metz, Lili Keren Szepel, Maarten Bosman, Maarten Sap, Maarten ter Hooft, Mahesh Paruchuri, Manal Feneq, Mantas Maravila, Marco Baranzan, Marco Marini, Marco Maru, Maria Jose Ramirez Quiroz, Marie Tilkka, Mario Giannelli, Martha Lewis, Martin Pottner, Matthew L. Lauvin, Matthias Hagen, Maylis Schubert, Medina Orduña Rautavaara, Melody Arnold, Melvin McElrath, Michael A. Yu, Michael Cohen, Michael Gu, Michael Beutelsky, Michael Starin, Michael Strube, Michael Smeydowski, Michele Bevilacqua, Michihito Yasunaga, Mihir Kale, Mike Cain, Mimer Xu, Ming Sargun, Mo Tward, Mohit Bansal, Moira Aminpour, Mor Geva, Mostafid Ghafar, Mukund Varma T, Nanyang Peng, Nathan Chi, Nayana Lee, Neta Gur-Ari Kradover, Nicholas Canova, Nicholas Roberts, Nick Doiron, Nikita Nangia, Niklas Deckert, Niklas Muenninghoff, Nirish Shrikir Kedar, Nivethita S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agbar, Omar Elbaghdadi, Omar Levy, Oran Etzioni, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Pang, Paul Fu Liang, Paul Vicol, Piyush Aljournabadi, Polyan Liao, Porey Liang, Peter Chang, Peter Eckersley, Phu Hoa Huu Phung Huang, Piotr Mikowski, Piyush Puri, Pooja Prasadipour, Priti Oli, Qianbo Mei, Qing Lyu, Qinfeng Chen, Robin Ranjale, Rachel Ema Radolph, Raef Gabriel, Rafael Haber, Ramon Elgo Elgado, Raphael Mollin, Rhytham Gang, Richard Barnes, RIT A. Samson, Riku Arakawa, Robbe Raymond, Robert Frank, Rohan Sikand, Roman Vukob, Roman Smolov, Roman LeRus, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, Ryan Chi, Ryan Lee, Ryan Soval, Ryan Tsohan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sajjan Aam, Sam Dillavos, Sam Shalifer, Sam Wiseman, Samuel Gasser, Samuel R. Bowman, Samuel S. Schoenholz, Saquyan Han, Sanjeev Kumar, Sarah A. Rose, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadoghi, Shadi Hashemi, Sharon Zhou, Shaohui Srivastava, Sherry Shi, Shikhar Singh, Shima Asadi, Shizhang Shao, Shi Shih Pichler, Shihua Tschalder, Shyam Upadhyay, Shyamkum (Shamun) Debnath, Siamak Shakeri, Simon Thormeyer, Susane Meli, Siva Reddy, Soha Pritika Makini, Soa-Hwan Lee, Spencer Torres, Sriharsha Warier, Stanislas Dehaene, Stefan Diele, Stefano Ermon, Stella Hideman, Stephanie Lin, Stephen Prasad, Steven T. Piantedosi, Stuart M. Shieber, Sumner Mishorghi, Svetlana Kirichenko, Swaroop Mishra, Tai Linson, Tai Schaefer, Tao Li, Tao Yu, Taryk Ali, Tama Hashimoto, To-Lin Wu, Tolo Enderoch, Theodore Rothchild, Thomas Phan, Tunde Wang, Tobias Minkoff, Timo Schick, Timofei Kornev, Timothy Tollu-Lawson, Tins Tansky, Tobias Gerstner, Treven Chang, Trishala Norraj, Tushar Khot, Tyler Shultz, Uri Shalun, Vidan Mera, Vera Dornberg, Verónica Nyman, Vikas Rastaki, Vinay Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikanth, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tang, Xinran Zhao, Xinyi Wu, Xudong Shen, Yafotah Yaghoobzadeh, Yair Laksy, Yanyan Qian, Yassam Bahet, Yefim Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hua, Yufang Hou, Yuntao Bai, Zachary Seld, Zheyan Zhao, Zijian Wang, Zijie J. Wang, Ziru Wang, Ziyi Wu.

Current Status

- 55 tasks created from 33 datasets
- Covering 6 dimensions
 - Content: speech command recognition
 - Speaker: speaker verification
 - Semantics: sarcasm detection
 - Degradation: noise SNR prediction
 - Paralinguistic: emotion recognition
 - Audio: environmental sound classification

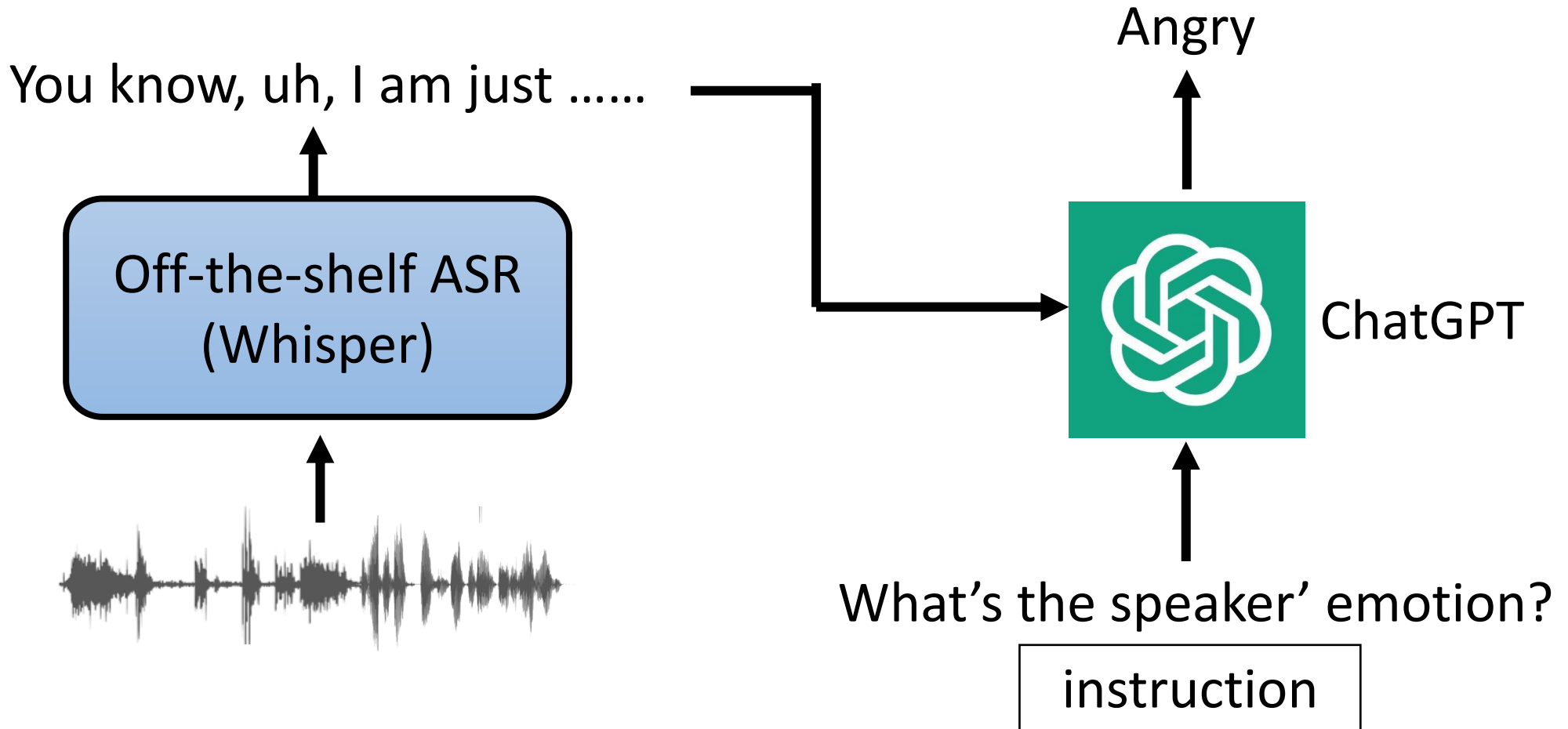
Everyone can add new tasks!

➡ **Dynamic**

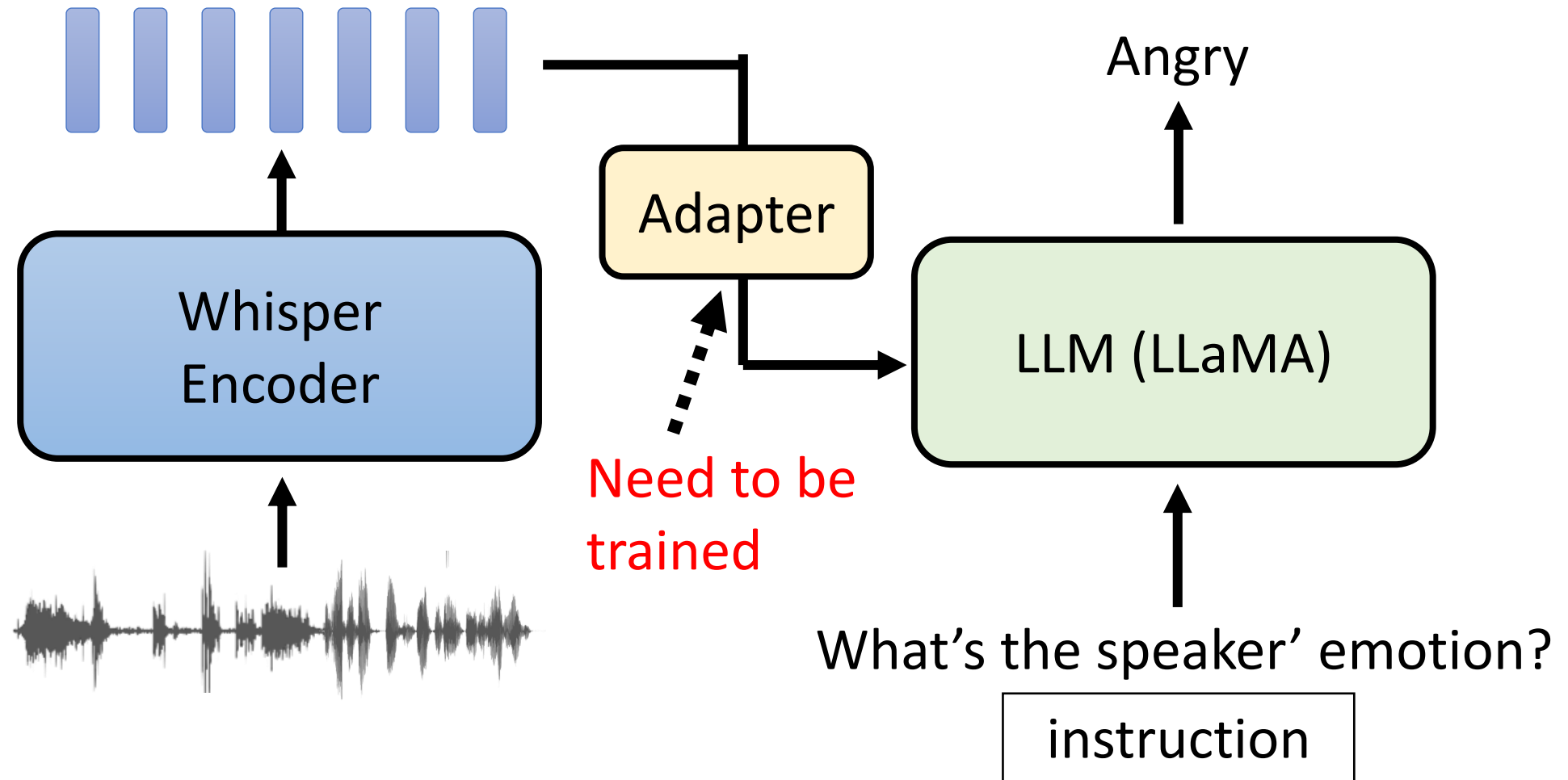


Let's work together!

Baseline: ASR + ChatGPT

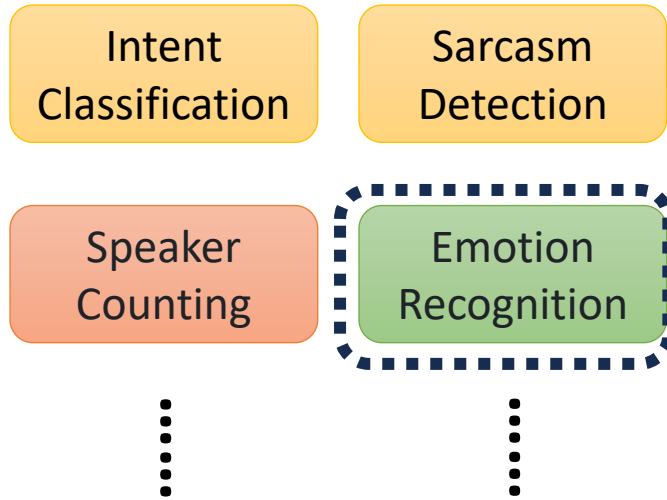


Baseline: Whisper Encoder + LLM



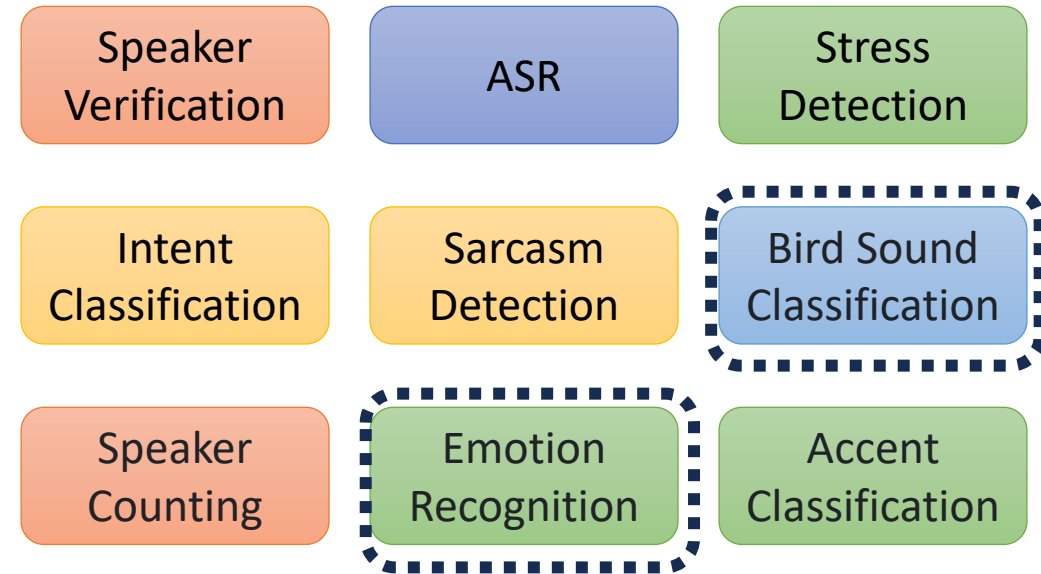
Training Tasks vs. Testing Tasks

Training Tasks (23 tasks)



Training model
parameters

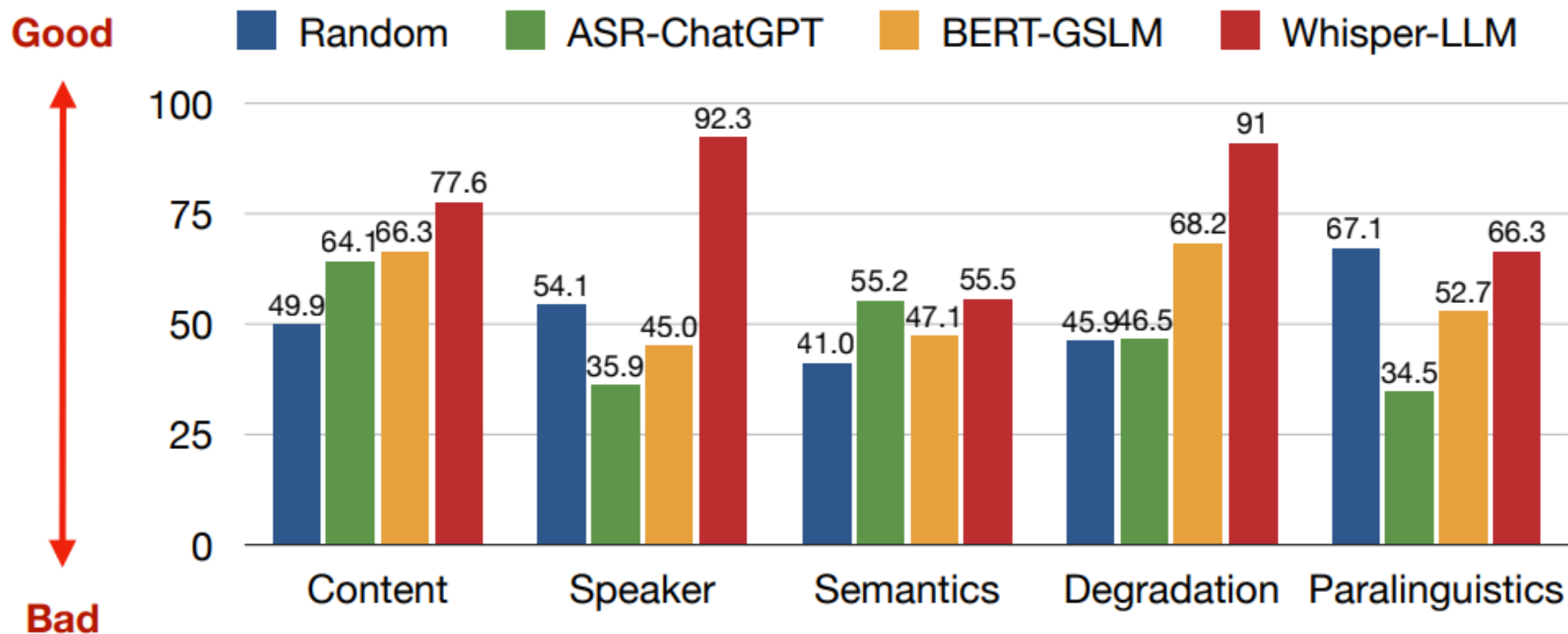
Dynamic-SUPERB (55 tasks)



Seen
(not the
same data)

Unseen

Overall Results



Dynamic-SUPERB

