# Spoken Conversational Agents with Large Language Models

**Chao-Han Huck Yang**[1]    **Andreas Stolcke**[2]    **Larry P Heck**[3]

NVIDIA Research[1]    Uniphore[2]    Georgia Institute of Technology[2]

hucky@nvidia.com    andreas.stolcke@uniphore.com    larryheck@gatech.edu

## 1   Introduction

Recent advancements in large language models (LLMs) with *voice interfaces* have garnered significant attention from both the research community and broader society. Closed-source models, such as GPT-4o and Gemini-1.5-pro, have demonstrated superior performance across classical speech tasks, including (i) speech recognition, (ii) translation, and (iii) spoken language understanding, significantly surpassing previous open-source benchmarks. Despite these strides, there remains a lack of comprehensive studies on the design and mechanisms underpinning the integration of speech modalities into LLMs for true multi-modal understanding. Additionally, the interaction between speech models and LLMs, particularly in the context of layered, self-play cognitive agents (Shah et al., 2018a,b), has only recently begun to be explored.

This tutorial aims to provide a thorough review of the historical trajectory of probabilistic language modeling (Jurafsky and Martin, 2024) for speech processing, offering insights that motivate the development of multi-agents system in conversational models. We will delve into advanced topics such as cross-modal adaptation, introducing on the theoretical foundations (Yang et al., 2021) required to align non-text modalities with textual representations. These concepts will be discussed alongside open-source, reproducible benchmarks (Chen et al., 2023a), providing a practical grounding for participants.

In addition, we will explore more recent trends toward end-to-end multi-modal speech-language models, emphasizing designs that utilize generative autoregressive approaches with joint speech-text tokenization. By examining both cascaded and end-to-end perspectives, this tutorial will equip participants with a comprehensive understanding of current strategies and open challenges in speech-augmented LLMs as shown in Figure 1.
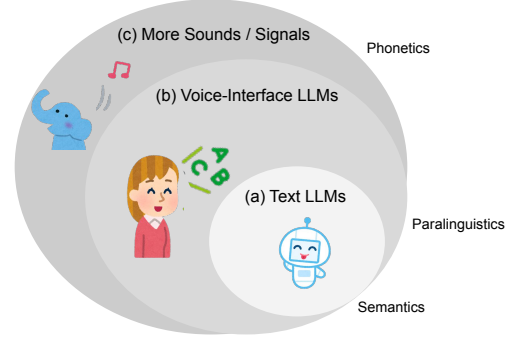


Figure 1: Examples of spoken conversational agent with different LLMs to understand different linguistic information from (a) semantics, (b) paralinguistics, and to (c) more phonetic signals.

## 2   Tutorial Outline

This three-hour tutorial will focus on the rapidly evolving field of speech-language modeling in the era of voice-interfacing LLMs and agent systems. Each core theme will be covered in a 35-minute segment, followed by a 10-minute Q&A and a 10-minute break. For each section, we will provide an overview of the relevant topics, followed by an in-depth exploration of key studies that shape the current landscape. The tutorial will conclude with a discussion of the open challenges and emerging research opportunities in this exciting and transformative area of joint speech-language modeling.

### 2.1   Language Modeling for Speech Processing Background, History, and Beyond

- **Probabilistic LMs for Speech Signals:** We will begin by discussing the early foundational work in Bayesian and n-gram-based language models, which were instrumental in building probabilistic frameworks for speech recognition tasks. These models laid the groundwork for current systems by modeling word sequences and handling uncertainties in speech inputs.

- **Contextual End-to-End Speech Models:** Next, we will explore the rise of end-to-end speech models that integrate contextual information directly into the modeling process. These models capture not only linguistic content but also paralinguistic features such as prosody, emotion, and speaker characteristics, allowing for more robust and nuanced speech understanding and generation.

- **Post-ASR LLM Correction and Fairness:** Finally, we will examine the role of LLMs in post-ASR error correction, where language models are used to refine and correct the outputs of traditional ASR systems. Special attention will be given to the fairness and bias issues that arise in speech processing, particularly with regard to speaker variability (e.g., accent, dialect, and sociolect). We will discuss how recent advances are addressing these challenges, and the ethical implications of ensuring equitable performance across diverse speech populations.

## 2.2 Large Language Models for Audio, Speech, and Conversational Signals

- **Theoretical Basics of LLMs Adaptation**: We will review some theoretical foundations of LLM adaptation and how these frameworks connect to LLMs with speech input. Topics include population risk measurement (Yang et al., 2021) and model transferability estimation (Chen et al., 2023b) from speech models to LLM adaptation, motivating different design pipelines (Radhakrishnan et al., 2023; Hu et al., 2024; Chen et al., 2024; Yang et al., 2023) of spoken agents.

- **Speech-Text Pre-training / Post-alignment**: Building on this, we will examine joint text-speech pre-training (Chiu et al., 2022; Barrault et al., 2023; Chen et al., 2022) methods, which have pushed the boundaries of multimodal understanding by combining speech and text learning objectives. The application of LLMs for voice quality estimation will also be discussed, showing how these models can assess and adapt to different speaker characteristics in real-time.

- **Multi-task Evaluation for Voice-LLMs:** Lastly, we will cover the latest advances in

joint generative translation models, which integrate both speech and text modalities for seamless, high-quality translation across languages. This section will provide a comprehensive look at how state-of-the-art voice-interfaced LLMs (Reid et al., 2024; Chu et al., 2023; Radford et al., 2023) are the processing and understanding of speech and conversational signals.

## 2.3 Spoken Dialogue Systems

- **Historical Foundations** The last decade saw significant progress in the creation and deployment of large-scale voice-powered AI virtual assistant technology, starting with Siri (SRI/Apple), and then followed by Cortana (Microsoft), Google Assistant, Alexa (Amazon), and Viv/Bixby (Samsung). This section of the tutorial will cover the technological innovations underlying this era of AI virtual assistants including intent detection and slot filling with RNNs (Mesnil et al., 2014), dialogue management with reinforcement learning (Shah et al., 2016), leveraging web-scale query-clicks (Hakkani-Tür et al., 2011; Tur et al., 2011), exploiting the semantic web (Heck and Hakkani-Tür, 2012), and incorporating knowledge graphs in spoken language understanding (Huang et al., 2015)

- **Current Trends** The current work in AI virtual assistants builds upon the voice-only systems of the last decade by leveraging LLMs to significantly improve the coverage and robustness of the spoken language understanding and dialogue state tracking components, in addition to substantial advancements in spoken language generation. This tutorial section provides an overview of existing LLMs and methodologies for adapting them to downstream tasks (Ni et al., 2021; Qin et al., 2023; Yi et al., 2024). It highlights recent advancements in multi-turn dialogue systems, encompassing both LLM-based open-domain dialogue (ODD) and task-oriented dialogue (TOD) systems, as well as relevant datasets and evaluation metrics. Additionally, it addresses emerging research challenges associated with the development of LLMs and the growing demands on multi-turn dialogue systems.

- **Future Directions** Future work will build on

the recent progress in LLMs and unify task-oriented and open-domain dialogue systems. This section of the tutorial will cover this technology shift and highlight some examples of promising future directions. Some examples include building new foundation LLMs pre-trained with conversational data (Jawale et al., 2024), e2e training of dialogues (Liu et al., 2017, 2018), and dialogue self-play (Shah et al., 2018a,b). Another significant future direction is situated dialogue systems: grounding LLM-based dialogue systems with content (web pages (Heck et al., 2013), lists (Bapna et al., 2017), forms (Heck et al., 2024), tables (Sundar and Heck, 2023; Sundar et al., 2024b,a), papers with figures, equations, tables (Sundar et al., 2024c), and retrieved documents (Reichman and Heck, 2024), vision (Hakkani-Tür et al., 2014), knowledge (Reichman et al., 2023), emotion (Reichman et al., 2024), and expression including facial/lip movement, facial/body expressions, and gestures (Punjwani and Heck, 2024b,a).

## 3 Tutorial Presenters

**Huck Yang** is a Senior Research Scientist at NVIDIA Research. He received his PhD degree from Georgia Institute of Technology. His research focuses on speech-language modeling, robust speech recognition, and multimodal post-editing models. He gave a tutorial at ASRU and Intetspeech 2023 on "Parameter-Efficient Language Modeling for Speech Processing," and a series of tutorials at ICASSP on "Prompt Learning for Speech-Language Models" from 2022 to 2024. He has served as area chairs in Interspeech, SLT, and ICASSP and worked full time at Amazon AGI.

**Andreas Stolcke** is a Distinguished AI Scientist and VP at Uniphore. He holds a PhD in Computer Science from UC Berkeley and has previously held positions at Amazon, SRI International, and Microsoft. His research spans language modeling, speech recognition, speaker identification, and speech translation. He is an IEEE Fellow and a Fellow of the International Speech Communication Association (ISCA), with long-term research interests in language modeling for speech processing, paralinguistics, and conversational AI.

**Larry Heck** is a Professor with a joint appointment in ECE and Interactive Computing, Chief

Scientist of Tech AI, and Executive Director of the Machine Learning Center at the Georgia Institute of Technology. He holds the Rhesa S. Farmer Distinguished Chair of Advanced Computing Concepts and is a Georgia Research Alliance Eminent Scholar. He received a PhD in Electrical Engineering from Georgia Institute of Technology, is an IEEE Fellow, and has previously held positions at SRI International, Nuance, Yahoo!, Microsoft, Google, Viv Labs, and Samsung. His career includes research on acoustics, active noise and vibration control, speaker recognition, web search, AI virtual assistants, NLP, and conversational AI.

## 4 Reading List and Prerequisite

This tutorial is designed for researchers and practitioners working at the intersection of conversational agents, and spoken information interactions. We assume attendees will have a foundational understanding of NLP and speech processing. While experience with system deployment is beneficial, we will provide a carefully paced introduction to core materials to accommodate a broader audience. Below are a few papers that offer important insights and foundations for this tutorial:

- Dialogue act modeling for automatic tagging and recognition of conversational speech (Stolcke et al., 2000)

- Is spoken language all-or-nothing? Implications for future speech-based human-machine interaction (Moore, 2017)

- Dialogue Learning with Human Teaching and Feedback in End-to-End Trainable Task-Oriented Dialogue Systems (Liu et al., 2018)

- Voice2series: Reprogramming acoustic models for time series classification (Yang et al., 2021)

- Hyporadise: An open baseline for generative speech recognition with large language models (Chen et al., 2023a)

- SpeechGPT: Empowering Large Language Models with Intrinsic Cross-Modal Conversational Abilities (Zhang et al., 2023)

**Breadth** While we will reference dozens of relevant papers throughout the tutorial, we plan to take a closer look at 7-8 key research papers in detail. Of these, only 1-2 will be directly authored

by the presenters, ensuring a broad and balanced exploration of the field.

# 5 Diversity Considerations

As we advance toward more conversational agents with LLMs, addressing diversity from spoken information is not just a matter of fairness, but of achieving the robustness and versatility necessary for real-world applications. Voice-interface based LLMs have made significant progress, but their limitations in handling diverse linguistic variations—accents, dialects, sociolects – present both ongoing challenges and opportunities for improvement.

- **Accents and Dialects: A Multidimensional Challenge**

  Speech processing models have traditionally struggled with accent and dialect diversity. While LLMs have shown impressive performance, they still reflect the biases of the datasets on which they are trained. These models often prioritize standard accents or high-resource languages, leading to suboptimal performance for speakers with less-represented accents. Moving forward, we need to systematically address this by designing models that generalize across diverse phonetic and prosodic patterns. Studies such as in-context learning based adaptation and retrieval-augmented clustering could play a crucial role in making LLMs more resilient to linguistic variation, allowing the models to dynamically adapt to previously unseen accents with minimal data.

- **Redefining Diversity-Oriented Evaluation Metrics for Spoken Conversational Agents**

  Standard evaluation metrics for LLMs in speech processing—such as word error rate (WER) for ASR—are often insufficient for capturing the true performance disparities across diverse linguistic groups. We must expand our evaluation frameworks to include fairness-driven metrics that assess model performance across different demographics, accents, and languages. Additionally, error analysis should go beyond quantitative measures and include qualitative insights into which user groups are more affected by model biases. Such nuanced evaluation will help guide the development of more equitable models.

- **Sociolects and Variability Across Social Dimensions for Conversational Agents**

  The diversity within a language itself—manifested through sociolects—poses an additional layer of complexity. Speech patterns vary based on age, gender, socio-economic background, and even profession. Traditional speech models tend to overfit to more homogenous data distributions, often reflecting the sociolects of dominant groups in the training datasets. Addressing this requires more than just adding diverse data; it demands sophisticated mechanisms like style transfer and speaker adaptation that allow models to process a wide spectrum of linguistic and paralinguistic cues. In particular, attention to prosody, intonation, and speaker emotion is key to better understanding socially or emotionally charged speech.

# 6 Ethics Statement

As we advance the integration of LLMs with spoken interactions, it is essential to proactively address the ethical implications associated with these developments. In this tutorial, we aim to raise awareness of these issues while equipping participants with strategies for responsible innovation. To reach a broad and diverse audience, we will promote the tutorial across various platforms. Our presenters include both *early-career* and *senior* researchers from both industry and academia. On speech privacy and security consideration, with the growing ubiquity of voice-interfaced systems in everyday life, both in public and private domains, protecting users' speech data from potential misuse is of critical importance. Voice data can contain sensitive personal information, and ensuring its secure handling is non-negotiable.

Our team brings substantial expertise in this area. For instance, Huck has served on the IEEE data collection committee, reviewing voice and signal data benchmarks with a focus on ethical data use. Andreas and Larry, both IEEE Fellows, have extensive experience in language modeling for speech processing, including speaker identification, which presents unique privacy challenges. Andreas gave a plenary talk on "*Speech-based and Multimodal Approaches for Human versus Computer Addressee Detection*" at EMNLP 2016.

# References

Ankur Bapna, Gokhan Tür, Dilek Hakkani-Tür, and Larry Heck. 2017. Towards zero-shot frame semantic parsing for domain scaling. In *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*, volume 2017, pages 2476–2480.

Loïc Barrault, Yu-An Chung, Mariano Coria Meglioli, David Dale, Ning Dong, Mark Duppenthaler, Paul-Ambroise Duquenne, Brian Ellis, Hady Elsahar, Justin Haaheim, et al. 2023. Seamless: Multilingual expressive and streaming speech translation. *arXiv preprint arXiv:2312.05187*.

Chen Chen, Yuchen Hu, Chao-Han Huck Yang, Sabato Marco Siniscalchi, Pin-Yu Chen, and Eng-Siong Chng. 2023a. Hyporadise: An open baseline for generative speech recognition with large language models. *Advances in Neural Information Processing Systems*, 36.

Chen Chen, Ruizhe Li, Yuchen Hu, Chao-han Huck Yang, Sabato Marco Siniscalchi, Pin-Yu Chen, and Ensiong Chng. 2024. It's never too late: Fusing acoustic information into large language models for automatic speech recognition. In *International Conference on Learning Representations*.

Zhehuai Chen, Yu Zhang, Andrew Rosenberg, Bhuvana Ramabhadran, Pedro Moreno, Ankur Bapna, and Heiga Zen. 2022. Maestro: Matched speech text representations through modality matching. *arXiv preprint arXiv:2204.03409*.

Zih-Ching Chen, Chao-Han Huck Yang, Bo Li, Yu Zhang, Nanxin Chen, Shuo-Yiin Chang, Rohit Prabhavalkar, Hung-yi Lee, and Tara N Sainath. 2023b. How to estimate model transferability of pre-trained speech models? *arXiv preprint arXiv:2306.01015*.

Chung-Cheng Chiu, James Qin, Yu Zhang, Jiahui Yu, and Yonghui Wu. 2022. Self-supervised learning with random-projection quantizer for speech recognition. In *International Conference on Machine Learning*, pages 3915–3924. PMLR.

Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. 2023. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. *arXiv preprint arXiv:2311.07919*.

Dilek Hakkani-Tür, Larry Heck, and Gokhan Tur. 2011. Exploiting query click logs for utterance domain detection in spoken language understanding. In *2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5636–5639. IEEE.

Dilek Hakkani-Tür, Malcolm Slaney, Asli Celikyilmaz, and Larry Heck. 2014. Eye gaze for spoken language understanding in multi-modal conversational interactions. In *Proceedings of the 16th International Conference on Multimodal Interaction*, pages 263–266.

Larry Heck and Dilek Hakkani-Tür. 2012. Exploiting the semantic web for unsupervised spoken language understanding. In *2012 IEEE Spoken Language Technology Workshop (SLT)*, pages 228–233. IEEE.

Larry Heck, Dilek Hakkani-Tür, Madhu Chinthakunta, Gokhan Tur, Rukmini Iyer, Partha Parthasarathy, Lisa Stifelman, Elizabeth Shriberg, and Ashley Fidler. 2013. Multi-modal conversational search and browse. In *CEUR Workshop Proceedings*, volume 1012, pages 96–101. CEUR-WS.

Larry Heck, Simon Heck, and Anirudh S Sundar. 2024. mforms: Multimodal form filling with question answering. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 11262–11271.

Yuchen Hu, Chen Chen, Chao-Han Yang, Ruizhe Li, Dong Zhang, Zhehuai Chen, and EngSiong Chng. 2024. GenTranslate: Large language models are generative multilingual speech and machine translators. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 74–90, Bangkok, Thailand. Association for Computational Linguistics.

Hongzhao Huang, Larry Heck, and Heng Ji. 2015. Leveraging deep neural networks and knowledge graphs for entity disambiguation. *arXiv preprint arXiv:1504.07678*.

Toshish Jawale, Chaitanya Animesh, Sekhar Vallath, Kartik Talamadupula, and Larry Heck. 2024. Are human conversations special? a large language model perspective. *arXiv preprint arXiv:2403.05045*.

Daniel Jurafsky and James H. Martin. 2024. *Speech and Language Processing (3rd edition draft)*. Available from https://web.stanford.edu/˜jurafsky/slp3/.

Bing Liu, Gokhan Tur, Dilek Hakkani-Tur, Pararth Shah, and Larry Heck. 2017. End-to-end optimization of task-oriented dialogue model with deep reinforcement learning. *arXiv preprint arXiv:1711.10712*.

Bing Liu, Gökhan Tür, Dilek Hakkani-Tur, Pararth Shah, and Larry Heck. 2018. Dialogue learning with human teaching and feedback in end-to-end trainable task-oriented dialogue systems. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2060–2069.

Grégoire Mesnil, Yann Dauphin, Kaisheng Yao, Yoshua Bengio, Li Deng, Dilek Hakkani-Tur, Xiaodong He, Larry Heck, Gokhan Tur, Dong Yu, et al. 2014. Using recurrent neural networks for slot filling in spoken language understanding. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 23(3):530–539.

Roger K Moore. 2017. Is spoken language all-or-nothing? implications for future speech-based human-machine interaction. *Dialogues with social robots: enablements, analyses, and evaluation*, pages 281–291.

Jinjie Ni, Tom Young, Vlad Pandelea, Fuzhao Xue, Vinay Vishnumurthy Adiga, and E. Cambria. 2021. Recent advances in deep learning based dialogue systems: a systematic survey. *Artificial Intelligence Review*, 56:3055–3155.

Saif Punjwani and Larry Heck. 2024a. Allo-ava: A large-scale multimodal conversational ai dataset for allocentric avatar gesture animation. *Preprint*, arXiv:2410.16503.

Saif Punjwani and Larry Heck. 2024b. Large body language models. *Preprint*, arXiv:2410.16533.

Libo Qin, Wenbo Pan, Qiguang Chen, Lizi Liao, Zhou Yu, Yue Zhang, Wanxiang Che, and Min Li. 2023. End-to-end task-oriented dialogue: A survey of tasks, methods, and future directions. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP*, pages 5925–5941, Singapore. Association for Computational Linguistics.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR.

Srijith Radhakrishnan, Chao-Han Huck Yang, Sumeer Khan, Rohit Kumar, Narsis Kiani, David Gomez-Cabrero, and Jesper Tegnér. 2023. Whispering llama: A cross-modal generative error correction framework for speech recognition. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10007–10016.

Benjamin Reichman and Larry Heck. 2024. Dense passage retrieval: Is it retrieving? In *Proceedings of the 2024 Empical Methods in Natural Language Processing (EMNLP 2024)*, Miami, Florida.

Benjamin Reichman, Kartik Talamadupula, Toshish Jawale, and Larry Heck. 2024. Reading with intent. *arXiv preprint arXiv:2408.11189*.

Benjamin Z Reichman, Anirudh Sundar, Christopher Richardson, Tamara Zubatiy, Prithwijit Chowdhury, Aaryan Shah, Jack Truxal, Micah Grimes, Dristi Shah, Woo Ju Chee, et al. 2023. Outside knowledge visual question answering version 2.0. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.

Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.

P. Shah, D. Hakkani-Tür, and L. Heck. 2016. Interactive reinforcement learning for task-oriented dialogue management. In *NIPS 2016 Deep Learning for Action and Interaction Workshop (Vol. 11)*.

Pararth Shah, Dilek Hakkani-Tür, Bing Liu, and Gokhan Tür. 2018a. Bootstrapping a neural conversational agent with dialogue self-play, crowdsourcing and on-line reinforcement learning. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 3 (Industry Papers)*, pages 41–51, New Orleans - Louisiana. Association for Computational Linguistics.

Pararth Shah, Dilek Hakkani-Tür, Gokhan Tür, Abhinav Rastogi, Ankur Bapna, Neha Nayak, and Larry Heck. 2018b. Building a conversational agent overnight with dialogue self-play. *Preprint*, arXiv:1801.04871.

Andreas Stolcke, Klaus Ries, Noah Coccaro, Elizabeth Shriberg, Rebecca Bates, Daniel Jurafsky, Paul Taylor, Rachel Martin, Carol Van Ess-Dykema, and Marie Meteer. 2000. Dialogue act modeling for automatic tagging and recognition of conversational speech. *Computational linguistics*, 26(3):339–373.

Anirudh Sundar, Christopher Richardson, William Gay, and Larry Heck. 2024a. itbls: A dataset of interactive conversations over tabular information. *arXiv preprint arXiv:2404.12580*.

Anirudh Sundar, Christopher Richardson, and Larry Heck. 2024b. gtbls: Generating tables from text by conditional question answering. *arXiv preprint arXiv:2403.14457*.

Anirudh Sundar, Jin Xu, William Gay, Christopher Richardson, and Larry Heck. 2024c. cpapers: A dataset of situated and multimodal interactive conversations in scientific papers. In *NeurIPS 2024*, Vancouver, Canada.

Anirudh S. Sundar and Larry Heck. 2023. cTBLS: Augmenting large language models with conversational tables. pages 59–70, Toronto, Canada.

Gokhan Tur, Dilek Hakkani-Tür, and Larry Heck. 2011. Sentence simplification for spoken language understanding. In *2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5628–5631. IEEE.

Chao-Han Huck Yang, Yile Gu, Yi-Chieh Liu, Shalini Ghosh, Ivan Bulyko, and Andreas Stolcke. 2023. Generative speech recognition error correction with large language models and task-activating prompting. In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 1–8. IEEE.

Chao-Han Huck Yang, Yun-Yun Tsai, and Pin-Yu Chen. 2021. Voice2series: Reprogramming acoustic models for time series classification. In *International conference on machine learning*, pages 11808–11819. PMLR.

Zihao Yi, Jiarui Ouyang, Yuwen Liu, Tianhao Liao, Zhe Xu, and Ying Shen. 2024. A survey on recent advances in llm-based multi-turn dialogue systems. *ArXiv*, abs/2402.18013.

Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. 2023. Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 15757–15773.