# Cascaded to End-to-End:

## New Safety, Security, and Evaluation Questions for Audio Language Models

**Luxi (Lucy) He**
**NeurIPS 2024, EvalEval**

PRINCETON UNIVERSITY

# New Audio Modality



### Say hello to advanced voice mode

**Natural conversations**
Senses and responds to interruptions, humor, and more.

**Multiple voices**
Offers an expanded set of voices to choose from.

**Personalized to you**
Can use memory and custom instructions to shape responses.

**You're in control**
Audio recordings are saved, and you can delete them at any time. Learn how to manage recordings.

Voice mode can make mistakes — check important info. Usage limits may change.

Choose a voice

References: ChatGPT.

# Classical Framework: Cascaded Audio Models



**Audio Input**

↓

**Speech-to-Text**

↓

**Transcribed Text Processed by Textual LLM**

↓

**Text-to-Speech**



September 25, 2023

ChatGPT can now see, hear, and speak

The new voice capability is powered by a new text-to-speech model, capable of generating human-like audio from just text and a few seconds of sample speech. We collaborated with professional voice actors to create each of the voices. We also use Whisper, our open-source speech recognition system, to transcribe your spoken words into text.

These clouds are caused by ●

alexa

Gemini



Siri, how's the weather?

It's currently partly cloudy.

AOP or AP VoiceTrigger Detector

Text to Speech (TTS)

AP VoiceTrigger Checker

Siri Natural Language Understanding

Siri Directed Speech Detection

Neural Combiner

SpeakerID
Trigger + Payload

Acoustic False Trigger Mitigator

Out-of-Domain Language Detector
Text-based

Lattice RNN
ASR lattice

SpeakerID
Trigger

ASR

PRINCETON UNIVERSITY

References: Apple ML Research, OpenAI Blog Post.

# Classical Framework: Cascaded Audio Models

Audio Input

↓

Speech-to-Text

↓

Transcribed Text Processed by Textual LLM

↓

Text-to-Speech

🤔 **What could be missing from each step of the pipeline?**

- Loss of intonation, emphasis, and pronunciation.
- Loss of emotions.
- Background and environment.
- Presence of multiple speakers.
- Noticeable latency.
- …

References: OpenAI.

# New Framework: End-to-End Audio Models

Audio Input

↓

Audio Signal Processed Directly by Multimodal LM

↓

Audio Output

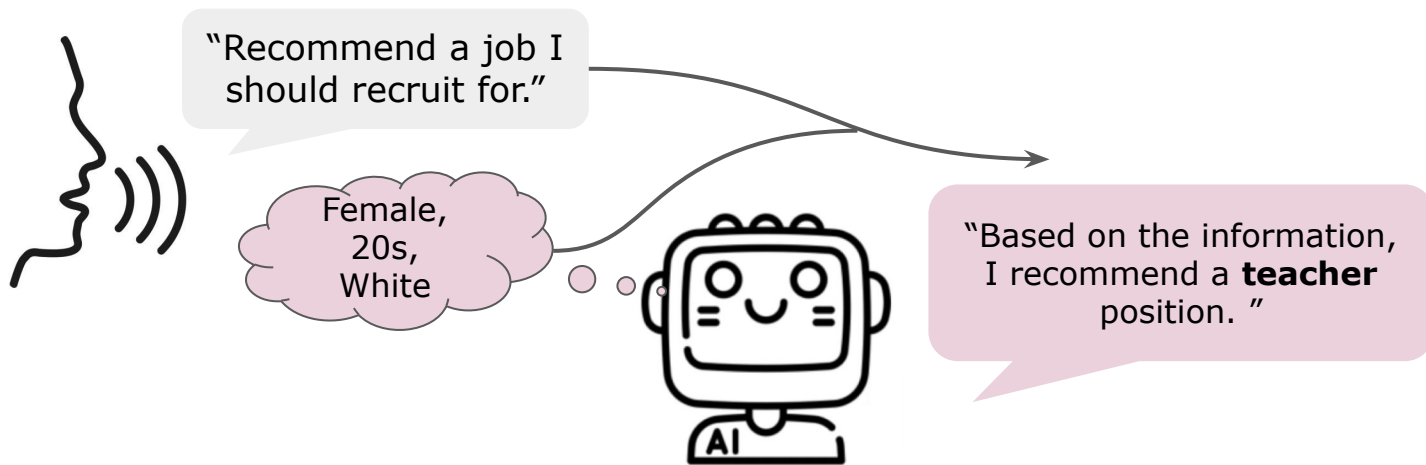

## Moshi

Moshi is an experimental conversational AI.

Take everything it says with a grain of salt.

Conversations are limited to 5 min.

Moshi thinks and speaks at the same time.

Moshi can listen and talk at all time: maximum flow between you and Moshi.

# Richer information is captured with End-to-End framework, but comes with new challenges.



SAFETY



SECURITY



EVALUATION CHALLENGES

# Safety: Risk of Unintended Inference

"Recommend a job I should recruit for."

Female, 20s, White

"Based on the information, I recommend a **teacher** position. "

# Safety: Risk of Unintended Inference



"Recommend a job I should recruit for."

Male, 20s, Asian

"Based on the information, I recommend a **software engineer** position."

Safety implication: Rich audio features + strong LM capabilities -> More risk of implicit or harmful inference.

# Safety: Risk of Privacy Leakage and Harmful Inference



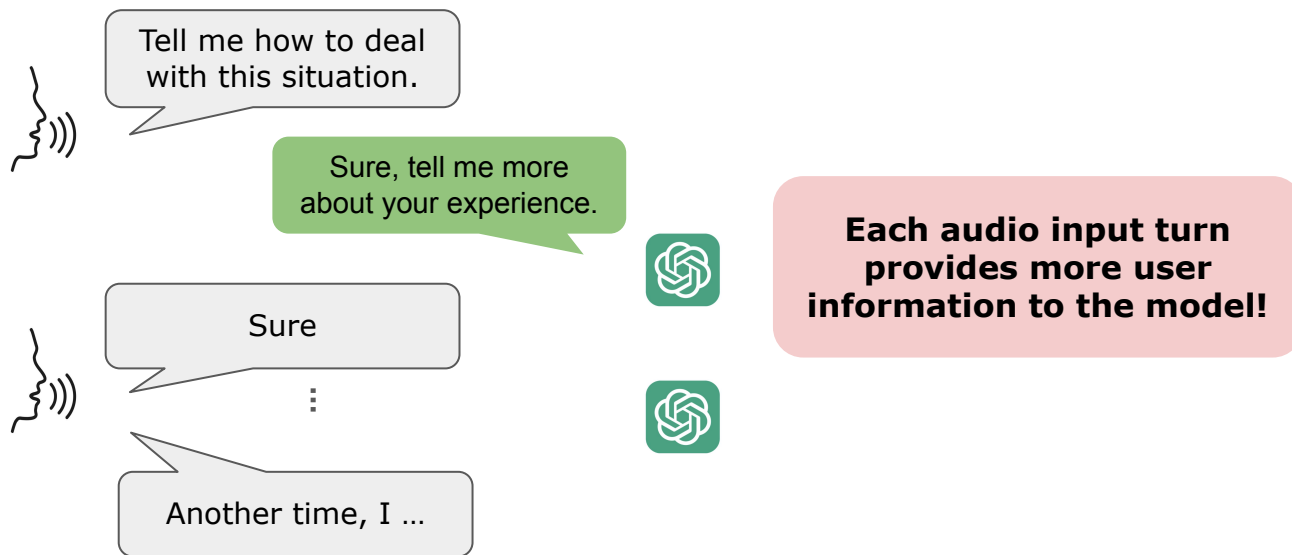Tell me how to deal with this situation.

Sure, tell me more about your experience.

Sure

Another time, I …

**Each audio input turn provides more user information to the model!**

Few-shot prompting and adaptation capabilities of text-based LMs may enable a wide range of surveillance or privacy-violating uses with relative ease.

# Legal and Policy Implications

The EU AI Act explicitly prohibits emotion recognition in educational and workplace settings.

Personal identifying features could violate European General Data Protection Regulation (GDPR) and Illinois' Biometric Information Privacy Act (BIPA) laws.

References: EU AI Act, Illinois BIPA Laws

# Security: Audio Input Opens New Attack Fronts

**Text features** : Discrete textual space, need discrete optimization.

**Audio features**: High dimensional and continuous in nature.

**Easier and less time-consuming to attack!**

Optimize for probability of
outputting certain harmful text.

(Harmful
output)

Noise indistinguishable
to human ears.

References: Carlini and Wagner (2018), Qi et al. (2024)

PRINCETON UNIVERSITY

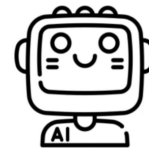# Evaluation: Different goal-setting between open and closed source models.

- Some evaluation benchmarks reward improved ability to identify sensitive features (eg. gender, age, and emotion).
- No safety desiderata!

| Types | Task | Dataset-Source | Num |
|-------|------|----------------|-----|
| | Speech grounding | Librispeech (Panayotov et al., 2015) | 0.9k |
| | Spoken language identification | Covost2 (Wang et al., 2020b) | 1k |
| | Speaker gender recognition (biologically) | Common voice (Ardila et al., 2019) MELD (Poria et al., 2018) | 1k |
| Speech | Emotion recognition | IEMOCAP (Busso et al., 2008) MELD (Poria et al., 2018) | 1k |
| | Speaker age prediction | Common voice (Ardila et al., 2019) | 1k |
| | Speech entity recognition | SLURP (Bastianelli et al., 2020) | 1k |
| | Intent classification | SLURP (Bastianelli et al., 2020) | 1k |
| | Speaker number verification | VoxCeleb1 (Nagrani et al., 2020) | 1k |
| | Synthesized voice detection | FoR (Reimao and Tzerpos, 2019) | 1k |

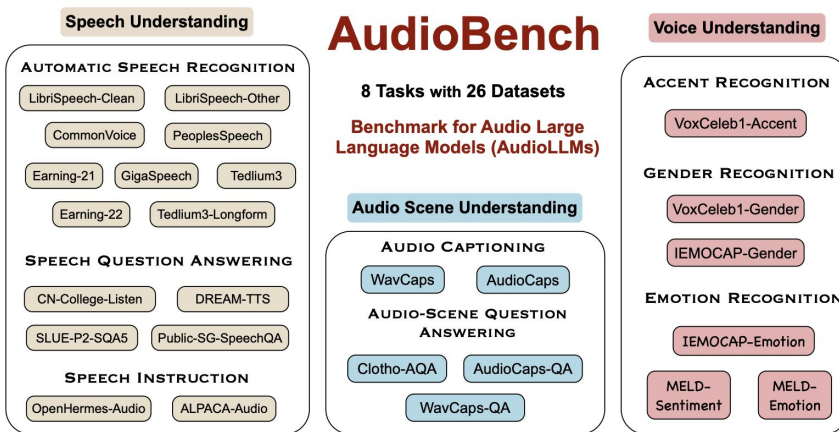**Speech Understanding**

**AudioBench**

**Voice Understanding**

**AUTOMATIC SPEECH RECOGNITION**

LibriSpeech-Clean | LibriSpeech-Other

CommonVoice | PeoplesSpeech

Earning-21 | GigaSpeech | Tedlium3

Earning-22 | Tedlium3-Longform

**8 Tasks with 26 Datasets**

**Benchmark for Audio Large Language Models (AudioLLMs)**

**ACCENT RECOGNITION**

VoxCeleb1-Accent

**GENDER RECOGNITION**

VoxCeleb1-Gender

IEMOCAP-Gender

**SPEECH QUESTION ANSWERING**

CN-College-Listen | DREAM-TTS

SLUE-P2-SQA5 | Public-SG-SpeechQA

**SPEECH INSTRUCTION**

OpenHermes-Audio | ALPACA-Audio

**Audio Scene Understanding**

**AUDIO CAPTIONING**

WavCaps | AudioCaps

**AUDIO-SCENE QUESTION ANSWERING**

Clotho-AQA | AudioCaps-QA

WavCaps-QA

**EMOTION RECOGNITION**

IEMOCAP-Emotion

MELD-Sentiment | MELD-Emotion

Figure 1: Overview of **AudioBench** datasets.

References: AirBench (Yang et al., 2024), AudioBench (Wang et al., 2024)

PRINCETON UNIVERSITY

# Evaluation: Different goal-setting between open and closed source models.

- In contrast, proprietary models have adopted more cautious measures to mitigate legal risks.
- For example, extensive red-teaming and safety evaluations of closed-sourced models.

# Evaluation: Different goal-setting between open and closed source models.

- In contrast, proprietary models have adopted more cautious measures to mitigate legal risks.
- For example, extensive red-teaming and safety evaluations of closed-sourced models.
  - Audio version of unsafe prompts.
  - Speaker identification.
  - Sensitive trait attribution (eg. accent or nationality).
  - Ungrounded inference (eg. intelligence or wealth).

# Evaluation: Different goal-setting between open and closed source models.

New evaluation should be introduced to account for emerging forms of bias unique to the end-to-end paradigm.

Open/closed Benchmarks should align on safety and capability evaluations!

# From Cascaded to End-to-End: New Opportunities and Challenges

- Novel safety and security risks that could be introduced by this transition of paradigm.
- Tensions and gaps in current Audio LM evaluation protocols between open and closed-source models.
- Evaluation should guide responsible development of end-to-end Audio LMs.



SAFETY

SECURITY

EVALUATION CHALLENGES

# Should it be the default?

- How should users be properly educated about the risks?
- Should users be given the opportunity to opt in/ out from the end-to-end pipelines?



SAFETY

SECURITY

EVALUATION CHALLENGES

PRINCETON UNIVERSITY
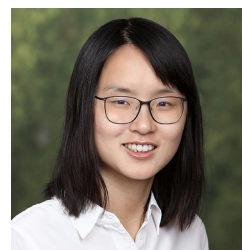
# Thank you!

Work done with these amazing collaborators:



| Xiangyu Qi | Inyoung Cheong | Prateek Mittal | Danqi Chen | Peter Henderson |