# Bridging the Perception Gap: A Lightweight Coarse-to-Fine Architecture for Edge Audio Systems

**Hengfan Zhang**[1]**, Yueqian Lin**[1]**, Hai "Helen" Li**[1]**, Yiran Chen**[1]

[1]**Duke University, Durham, NC, USA**

arXiv:2601.15676v1 [cs.SD] 22 Jan 2026

## Abstract

Deploying Audio-Language Models (Audio-LLMs) on edge infrastructure exposes a persistent tension between *perception depth* and *computational efficiency*. Lightweight local models tend to produce "passive perception"—generic summaries that miss the subtle evidence required for multi-step audio reasoning—while indiscriminate cloud offloading incurs unacceptable latency, bandwidth cost, and privacy risk. We propose **CoFi-Agent** (Tool-Augmented Coarse-to-Fine Agent), a hybrid architecture targeting edge servers and gateways. It performs fast local perception and triggers *conditional* forensic refinement only when uncertainty is detected. CoFi-Agent runs an initial single-pass on a local 7B Audio-LLM, then a cloud controller gates difficult cases and issues lightweight plans for on-device tools such as temporal re-listening and local ASR. On the **MMAR** benchmark, CoFi-Agent improves accuracy from 27.20% to **53.60%**, while achieving a better accuracy–efficiency trade-off than an always-on investigation pipeline. Overall, CoFi-Agent bridges the perception gap via tool-enabled, conditional edge–cloud collaboration under practical system constraints.

## 1. Introduction

Edge audio systems—from on-premise security gateways and smart-home hubs to autonomous service robots—must operate under strict system constraints. Latency, bandwidth, and acoustic privacy are first-order requirements for real-time applications, ranging from interactive voice assistants [1] to on-premise security gateways. Meanwhile, user queries increasingly demand *reasoning* over complex acoustic scenes: long recordings, heavy background noise, overlapping speakers, and rare but decisive events. Recent edge-AI surveys also highlight the growing tension between advanced foundation models and the practical limits of resource-constrained deployments [2].

Most Audio-Language Models (Audio-LLMs) still follow a *single-pass* paradigm: the model encodes the waveform once and generates a static explanation or answer. This design is computationally predictable, but it frequently fails on reasoning-oriented queries that require fine-grained evidence. When the decisive cue is weak (e.g., a short utterance) or temporally localized (e.g., event ordering), single-pass models cannot "go back" to verify the signal. As a result, they may produce plausible-but-unsupported rationales, creating a persistent *perception gap* between what the audio contains and what the system claims.

A naive fix is to always offload to the cloud or to always execute an expensive tool pipeline (e.g., ASR + multi-pass re-listening). However, always-on strategies are misaligned with edge constraints: they increase average latency and bandwidth usage, and expand the exposure surface of sensitive acoustic content. Moreover, tool execution can inject noise into easy samples (e.g., spurious ASR text in near-silence), which may even *hurt* accuracy when the baseline is already correct. Prior work in pervasive/edge sensing also suggests that pushing

computation and interaction *toward the point of collection* can be beneficial for practical deployment, while still requiring careful system design to avoid unnecessary overheads [3].

We propose **CoFi-Agent** (Tool-Augmented Coarse-to-Fine Agent), a hybrid architecture designed to close the perception gap *without* paying the always-on cost. CoFi-Agent performs a fast local perception pass for every query and *conditionally* escalates only uncertain cases for targeted refinement. During refinement, a cloud controller emits lightweight investigation plans, and the edge node executes on-device tools such as temporal re-listening and local ASR to extract high-value evidence. Crucially, raw audio remains on-premise; only compact evidence (e.g., transcripts and tool summaries) is transmitted for cloud reasoning, thereby prioritizing *acoustic privacy*.

**Contributions.** This paper makes three contributions:

- **System architecture:** a local-first coarse-to-fine edge–cloud workflow for audio reasoning, targeting high-performance edge gateways.

- **Tool-enabled refinement:** targeted re-listening and on-device ASR to recover query-relevant evidence while keeping raw audio local.

- **Benchmark evidence:** on MMAR (N=1,000), CoFi-Agent improves accuracy from 27.20% to 53.60% and achieves a stronger accuracy–efficiency trade-off than an always-on investigation baseline.

We quantify this efficiency advantage via an accuracy–latency trade-off analysis (Fig. 1) and present the overall system workflow in Fig. 2.

## 2. Related Work
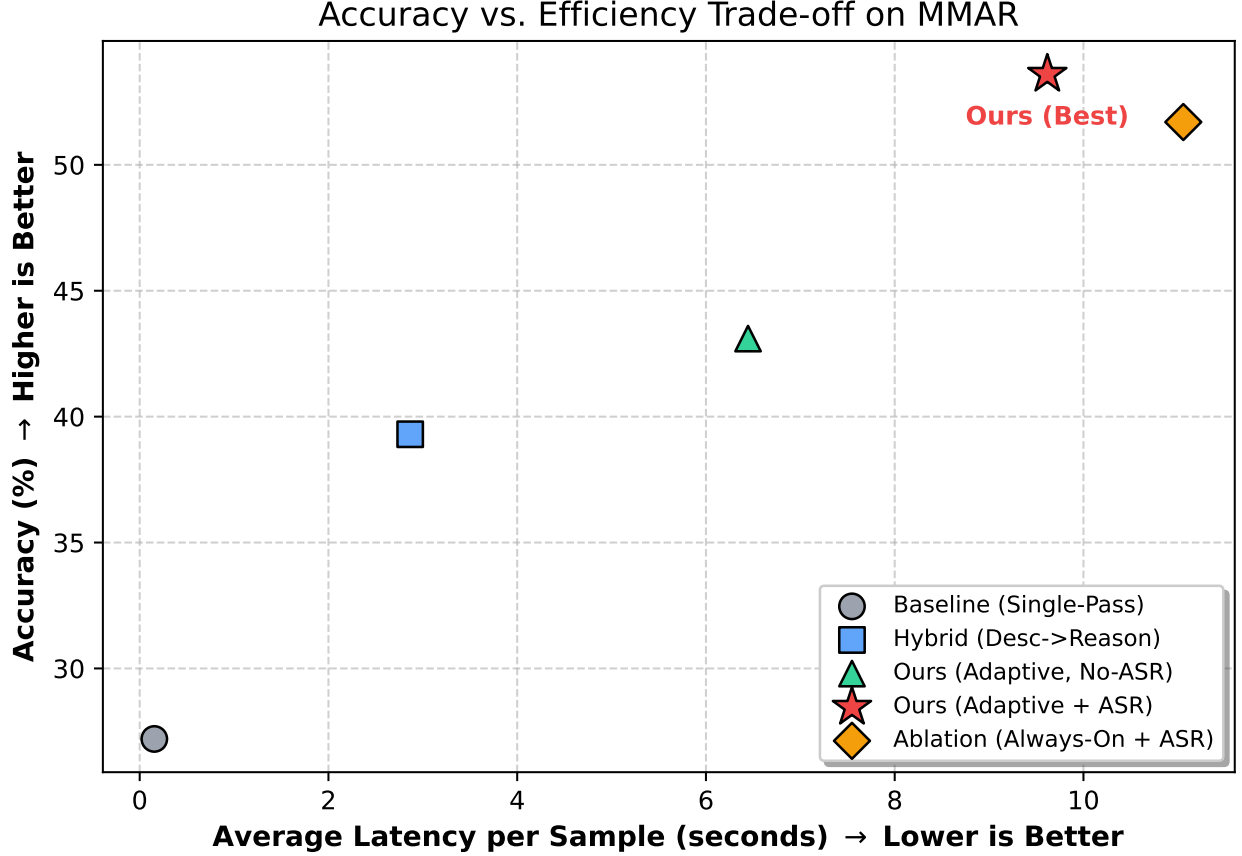
### 2.1. Audio-Language Models(Audio-LLMs)

The integration of audio encoders with Large Language Models has enabled open-ended acoustic understanding. Models like Qwen-Audio [4], LTU [5], and SALMONN [6], and recent multimodal memory architectures [7] map audio features into the LLM's embedding space, allowing for instruction following. However, these models predominantly operate in a *single-pass* manner. While efficient for simple captioning, they lack the iterative reasoning capabilities required for complex forensic tasks [8], where evidence is often subtle or temporally scattered.

### 2.2. Agentiv AI and Tool Use

Recent works have empowered LLMs with external tools (ASR, Python, Search) [9, 10]. Systems like HuggingGPT [11] and AudioGPT [12] orchestrate specialized models to solve complex queries. Yet, most agentic frameworks assume a cloud-native environment with unlimited bandwidth and compute. They often invoke tools indiscriminately, which, in an edge context, leads to prohibitive latency and privacy violations. CoFi-Agent differentiates itself by being *conditional* and *local-first*, triggering tool usage only when the lightweight perception model expresses uncertainty.

### 2.3. Efficient Edge Intelligence

Traditional efforts to deploy deep learning on the edge have focused largely on *static model compression*. Techniques such as quantization (e.g., LLM.int8() [13], AWQ [14]), pruning [15, 16], and knowledge distillation [17] are effective at reducing the memory footprint of individual models. However, these methods suffer from a

**Fig. 1**: Accuracy–efficiency trade-off on MMAR (N=1,000). Adaptive gating achieves higher accuracy and lower average latency than always-on investigation.

fundamental inefficiency: they apply the same computational budget to every input, regardless of difficulty. This results in *semantic redundancy*, where simple queries (e.g., clear speech) are processed with the same heavy parameters as complex forensic tasks. Our work addresses this by shifting focus to *dynamic inference*. Aligning with *Early-Exit* architectures (e.g., BranchyNet [18]) and classic *Cascade Systems* [19], CoFi-Agent implements a "semantic cascade": a lightweight perceptual check that exits early for the majority of easy samples, triggering expensive reasoning tools only for the "hard" tail of the distribution.
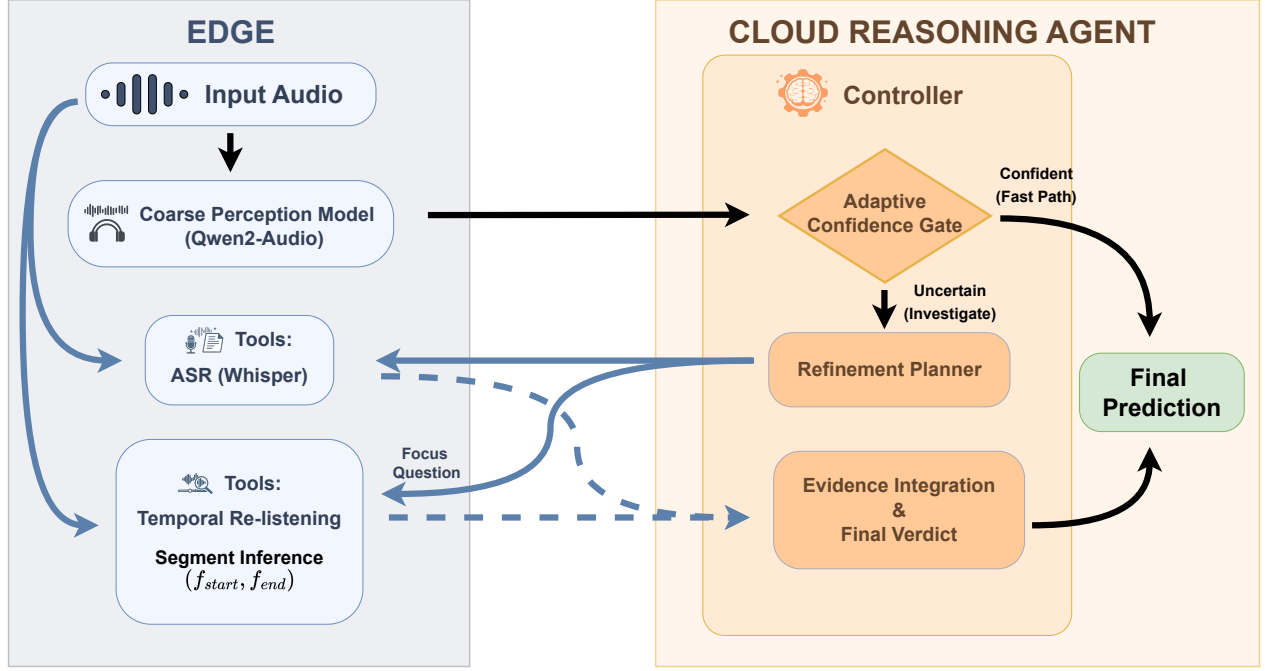
## 3. Methodology

We consider audio reasoning as predicting $y \in \mathcal{Y}$ given audio $A$ and query $Q$ while minimizing end-to-end cost $C$ (latency). CoFi-Agent decomposes inference into coarse perception and conditional refinement.

### 3.1. Stage 0: Edge Coarse Perception

A compact on-device Audio-LLM $\mathcal{M}_{\text{edge}}$ performs single-pass inference:

$$(s_0, p_0) = \mathcal{M}_{\text{edge}}(A, Q), \tag{1}$$

**Fig. 2**: CoFi-Agent overview. A local coarse perception model answers easy queries via a Fast Path. A cloud confidence gate escalates uncertain cases and emits lightweight refinement plans for on-device tools (temporal re-listening and ASR). Raw audio remains on-device; only compact evidence (e.g., transcripts, tool summaries) is shared for cloud reasoning.

where $p_0$ is the initial answer and $s_0$ is a short rationale/summary.

## 3.2. Adaptive Confidence Gate

A cloud controller evaluates ambiguity and self-consistency:

$$u = \mathcal{G}(s_0, Q, p_0), \qquad u \in \{0, 1\}. \tag{2}$$

If $u = 0$, return $y = p_0$ (Fast Path). If $u = 1$, trigger refinement (Investigate Path).

**Implementation of $\mathcal{G}$.** We implement $\mathcal{G}$ as a lightweight prompt-based classifier. It checks for uncertainty cues (hedging), missing evidence, and logical inconsistencies. In our experiments, $\mathcal{G}$ escalates approximately 62% of MMAR samples to the investigation path. Qualitative analysis shows that *false escalations* typically occur on low-SNR non-speech clips where the baseline is correct but "hedges" its language, while *missed escalations* occur when the baseline hallucinates confidently on speech-heavy questions.

## 3.3. Stage 1: Cloud-Guided Refinement Planning

To enable temporal refinement without uploading waveforms, CoFi-Agent uses an on-device segment proposer $\mathcal{P}$. We implement $\mathcal{P}$ with a hybrid strategy: (i) energy-based segmentation for distinct events, and (ii) uniform sliding windows ($K = 4$, 3s duration) for short clips ($< 12s$). For clips longer than 12s, windows are placed at relative percentiles (10%, 30%, 50%, 70%) to ensure global coverage. The cloud receives only **region-of-interest (ROI)** metadata (segment timestamps) and selects an ROI index $i^*$ and a focused sub-query $q_{\text{focus}}$.

### 3.4. Stage 2: On-Device Tool Execution

#### 3.4.1. Tool: Temporal Re-listening

The edge device runs targeted inference on the selected segment:

$$e_{\text{audio}} = \mathcal{M}_{\text{edge}}(A[t_s^{(i^*)} : t_e^{(i^*)}], q_{\text{focus}}). \tag{3}$$

#### 3.4.2. Tool: On-Device ASR (Whisper)

For speech-heavy queries, the device runs local ASR:

$$T_{\text{text}} = \text{ASR}_{\text{local}}(A \text{ or } A[t_s^{(i^*)} : t_e^{(i^*)}]). \tag{4}$$

### 3.5. Evidence Integration and Verdict

A cloud reasoner combines $(s_0, e_{\text{audio}}, T_{\text{text}}, Q)$ to output:

$$y_{\text{final}} = \mathcal{M}_{\text{cloud}}(s_0, e_{\text{audio}}, T_{\text{text}}, Q). \tag{5}$$

### 3.6. Privacy and Bandwidth Efficiency

Unlike cloud-native approaches that continuously stream raw waveforms, CoFi-Agent adheres to the principle of *Data Minimization*.

- **Acoustic Isolation:** Biometric markers (e.g., voiceprints) and irrelevant background environmental sounds never leave the local device.

- **Symbolic Transmission:** The cloud receives only compact, symbolic representations ($T_{\text{text}}$). We acknowledge that transcripts may still contain sensitive semantic information; however, this text-only format allows for efficient on-device PII redaction before upload, which is computationally infeasible for raw audio.

## 4. Experimental Evaluation

### 4.1. Benchmark and Metrics

We evaluate on the **MMAR** benchmark [20] (**N=1,000**). Each sample contains an audio clip, a natural-language question, and multiple-choice candidates. We report (i) **Accuracy** and (ii) **Cost** as the average end-to-end wall-clock latency (seconds/sample, Batch Size=1), including local inference, optional tool execution, network overhead, and cloud reasoning.

### 4.2. System Setup and Implementation Details

**Edge Environment Emulation.** We emulate a high-performance edge gateway using an **NVIDIA Quadro RTX 6000** (24GB VRAM). Each CoFi-Agent instance is allocated a single GPU to reflect an "edge-node budget." All models are kept warm in memory to avoid loading overhead. The local perception backbone is **Qwen2-Audio-7B-Instruct** [4] in FP16.
**On-Device Tools.** We use **Whisper-small** [21] as the local ASR tool, executed on the same edge GPU. Temporal re-listening re-invokes the local Audio-LLM on a selected ROI.

**Table 1**: Accuracy and Cost on MMAR (N=1,000).

| Method | Acc. | Cost |
|---|---|---|
| Qwen2-Audio (Baseline) | 27.20% | 0.155s |
| Hybrid (Describe→Reason) | 39.30% | 2.866s |
| CoFi-Agent (Adaptive Re-listen) | 43.10% | 6.446s |
| CoFi-Agent (Always-On + ASR) | 51.70% | 11.058s |
| **CoFi-Agent (Adaptive + ASR)** | **53.60%** | **9.617s** |

**Cloud Controller.** The cloud-side controller (gating, planning, and final verdict) uses **OpenAI GPT-4o** [22] with Temperature=0 for reproducibility. Network RTT from the edge node to the US-East region is approx. 15ms (p50) and 45ms (p95).

## 4.3. Main Results

Table 1 summarizes the overall accuracy and end-to-end latency on MMAR (N=1,000). CoFi-Agent (Adaptive + ASR) achieves the best accuracy (53.60%) while also improving the accuracy–latency operating point compared to the always-on investigation variant. This trend is consistent with the trade-off curve shown in Fig. 1.
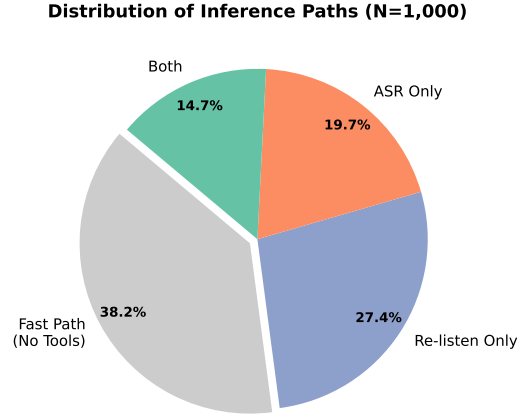
## 4.4. Why Adaptive Beats Always-On

Always-on investigation increases coverage but can inject tool noise into easy samples (e.g., ASR hallucinations in silence), and it pays the tool cost even when the baseline is already correct. Adaptive gating escalates only uncertain cases (about 62% on MMAR), yielding a better balance of accuracy and average latency. Table 2 reports the latency breakdown for CoFi-Agent (Adaptive + ASR).

**Table 2**: Latency Breakdown (CoFi-Agent Adaptive + ASR).

| Component | Latency (s) |
|---|---|
| Stage 0 (Edge Perception) | 0.16 |
| Network (RTT + Upload) | 0.20 |
| Cloud Gate (Classification) | 0.60 |
| On-Device Tools (Whisper/Re-listen) | 1.85 |
| Cloud Reasoning (Generation) | 6.81 |
| **Total** | **9.62** |

## 4.5. Tool-Usage Distribution Under Adaptive Gating

To better understand where the compute is spent, we measure the distribution of inference paths under adaptive gating. As shown in Fig. 3, **38.2%** of samples finish in the Fast Path with **no tools**, while **61.8%** are escalated. Among escalated samples, temporal re-listening alone accounts for **27.4%**, ASR alone for **19.7%**, and invoking *both* tools for **14.7%**. This suggests that (i) a substantial portion of queries are easy enough for single-pass local perception, and (ii) when escalation is needed, many cases can be resolved with a single lightweight tool rather than a full always-on pipeline.

**Fig. 3**: Distribution of inference paths under CoFi-Agent adaptive gating on MMAR (N=1,000).

## 4.6.  Impact of ASR Tools

Adding on-device ASR yields a +10.50% absolute gain (43.10% → 53.60%) compared to adaptive re-listening alone (Table 1), confirming that resolving speech semantics is decisive for a large fraction of MMAR questions. At the same time, Fig. 3 shows ASR is not universally required: adaptive routing avoids paying ASR cost on samples that do not benefit from it.

## 5.  Case Studies

We analyze three specific samples from the MMAR dataset to demonstrate how CoFi-Agent corrects perception errors using tool-augmented refinement.

## 5.1.  Sample A (Conversational Implication)

For the query, "Do the first man's two statements to the other person have the same implied meaning?", the single-pass baseline incorrectly predicts "Same." The model likely aggregates the general friendly tone of the voice, failing to distinguish the semantic shift between the two farewells. CoFi-Agent triggers the ASR tool, which reveals the exact phrasing: the first utterance is a standard "Have a good day," while the second is the ominous "Enjoy the next 24 hours of your life." By accessing this textual evidence, the cloud reasoner detects the shift in meaning and correctly predicts "Different."

## 5.2.  Sample B (Keyword Verification)

When asked, "Does the Chinese the speaker wants to show to their father include the word 'Ni Hao'?", the baseline model incorrectly asserts that the word is present. This error likely stems from the high prior probability of common greetings appearing in the context of a foreigner speaking Chinese. The investigation path executes a negative verification using the ASR transcript: *"...Dad, I'm in China, man... Erho? Erho. Okay, good..."*. The cloud reasoner confirms the specific target word is absent from the text and corrects the verdict to "Does not include."

### 5.3. Sample C (Semantic Trick/Wifi Password)

A query asking "What is the wifi password?" exposes a failure in parsing semantic riddles. The baseline model interprets the barista's response—"You have to buy Smoothie first"—as a refusal or a precondition. However, the refinement step captures the full dialogue via ASR: *"Can I get the Wi-Fi password please? You have to buy Smoothie first... Now can I get the password? You have to buy Smoothie first. I just did. That's the best word."* The cloud reasoner parses the humor in the transcript and correctly extracts the phrase "Youhavetobuysmooziefirst" as the password string.

### 5.4. Failure Mode Analysis

We observe three common failure cases that explain most remaining errors. **(1) ASR under extreme noise.** Under very low SNR conditions (< 0 dB), Whisper may produce corrupted transcripts, which misleads the cloud reasoner. **(2) Missed short events in long recordings.** The heuristic segment proposer may skip brief but decisive events in long clips (> 1 min), causing re-listening to focus on irrelevant regions. **(3) Knowledge-mismatch queries.** A subset of queries require external knowledge not present in the audio; even perfect perception cannot resolve these cases.

## 6. Conclusion

CoFi-Agent addresses the tension between perception depth and efficiency in edge audio systems by introducing a coarse-to-fine architecture that keeps raw waveforms on-premise and invokes cloud reasoning only for high-entropy queries (approx. 62% of cases). On the MMAR benchmark, this approach nearly doubles the accuracy of a local 7B model (27.20% $\rightarrow$ 53.60%) while maintaining a viable latency profile. Future work will focus on minimizing decision overhead via learnable gating, jointly training the segment proposer with the reasoning module, and extending this paradigm to bandwidth-constrained Edge Video scenarios. Ultimately, CoFi-Agent demonstrates that the future of Edge AI relies on adaptive system design—knowing when to think fast locally, and when to think slow via the cloud.

## References

[1] Yueqian Lin, Zhengmian Hu, Jayakumar Subramanian, Qinsi Wang, Nikos Vlassis, Hai Li, and Yiran Chen. Asyncvoice agent: Real-time explanation for llm planning and reasoning. In *IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2025. Demo Track.

[2] Mozhgan Navardi, Romina Aalishah, Yuzhe Fu, Yueqian Lin, Hai Li, Yiran Chen, and Tinoosh Mohsenin. Genai at the edge: Comprehensive survey on empowering edge devices. In *AAAI Spring Symposium Series, GenAI@Edge Workshop*, 2025.

[3] K. Woodward, E. Kanjo, and A. Oikonomou. LabelSens: Enabling Real-time Sensor Data Labelling at the point of Collection on Edge Computing. *arXiv preprint arXiv:1910.01400*, 2019.

[4] Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen2-audio technical report. 2024.

[5] Y. Gong et al. Listen, Think, and Understand. In *International Conference on Learning Representations (ICLR)*, 2024.

[6] C. Tang et al. SALMONN: Towards Generic Hearing Abilities for Large Language Models. In *International Conference on Learning Representations (ICLR)*, 2024.

[7] Yueqian Lin, Qinsi Wang, Hancheng Ye, Yuzhe Fu, Hai Li, and Yiran Chen. Hippomm: Hippocampal-inspired multimodal memory for long audiovisual event understanding. *arXiv preprint arXiv:2504.10739*, 2025.

[8] Yueqian Lin, Zhengmian Hu, Qinsi Wang, Yudong Liu, Hengfan Zhang, Jayakumar Subramanian, Nikos Vlassis, Hai Li, and Yiran Chen. Voice evaluation of reasoning ability: Diagnosing the modality-induced performance gap. *arXiv preprint arXiv:2509.26542*, 2025.

[9] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. Re-Act: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.

[10] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.

[11] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving AI tasks with ChatGPT and its friends in hugging face. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.

[12] Rongjie Huang, Mingze Li, Dongchao Yang, Jiatong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu, Zhiqing Hong, Jiawei Huang, Jinglin Liu, Yi Ren, Yuexian Zou, Zhou Zhao, and Shinji Watanabe. Audiogpt: Understanding and generating speech, music, sound, and talking head. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 23802–23804, 2024. doi: 10.1609/AAAI.V38I21.30570.

[13] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. LLM.int8(): 8-bit matrix multiplication for transformers at scale. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.

[14] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. AWQ: Activation-aware weight quantization for on-device LLM compression and acceleration. In *Proceedings of Machine Learning and Systems*, volume 6, 2024.

[15] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. In *International Conference on Learning Representations (ICLR)*, 2016.

[16] Yueqian Lin, Yuzhe Fu, Jingyang Zhang, Yudong Liu, Jianing Sun, Hai Li, and Yiran Chen. Speechprune: Context-aware token pruning for speech information retrieval. In *IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2025.

[17] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015.

[18] Surat Teerapittayanon, Bradley McDanel, and H. T. Kung. Branchynet: Fast inference via early exiting from deep neural networks. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pages 2464–2469, 2016. doi: 10.1109/ICPR.2016.7900006.

[19] Paul Viola and Michael Jones. Rapid object detection using a boosted cascade of simple features. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages I–511–I–518, 2001.

[20] Z. Ma et al. MMAR: A Challenging Benchmark for Deep Reasoning in Speech, Audio, Music, and Their Mix. In *Advances in Neural Information Processing Systems (NeurIPS), Datasets and Benchmarks Track*, 2025.

[21] A. Radford et al. Robust Speech Recognition via Large-Scale Weak Supervision. In *Proceedings of the 40th International Conference on Machine Learning (ICML)*, volume 202 of *Proc. Mach. Learn. Res.*, pages 28492–28518, 2023.

[22] OpenAI. Hello GPT-4o. https://openai.com/index/hello-gpt-4o/, May 2024. Accessed: 2025-12-22.