

# Speech-Hands: A Self-Reflection Voice Agentic Approach to Speech Recognition and Audio Reasoning with Omni Perception

Anonymous ACL submission

## Abstract

We introduce a voice-agentic framework that learns one critical omni-understanding skill: knowing when to trust itself versus when to consult external audio perception. Our work is motivated by a crucial yet counterintuitive finding: naively fine-tuning an omni-model on both speech recognition and external sound understanding tasks often degrades performance, as the model can be easily misled by noisy hypotheses. To address this, our framework, Speech-Hands, recasts the problem as an explicit self-reflection decision. This learnable reflection primitive proves effective in preventing the model from being derailed by flawed external candidates. We show that this agentic action mechanism generalizes naturally from speech recognition to complex, multiple-choice audio reasoning. Across the OpenASR leaderboard, Speech-Hands consistently outperforms strong baselines by 12.1% WER on seven benchmarks. The model also achieves 77.37% accuracy and high F1 on audio QA decisions, showing robust generalization and reliability across diverse audio question answering datasets. By unifying perception and decision-making, our work offers a practical path toward more reliable and resilient audio intelligence.

## 1 Introduction

Omni-modal models (Xie and Wu, 2024; OpenAI et al., 2024; Xu et al., 2025a; Li et al., 2025b) that jointly process audio and text have unified a range of audio understanding tasks, including automatic speech recognition (ASR), temporal sound event reasoning, and knowledge-heavy question answering. However, human perception is not naturally perfect at understanding acoustic patterns across different resolutions at soundscape (Calcus, 2024). For example, while professional speech interpreters could produce superior ASR transcriptions, this specialized ability does not guarantee a comparable aptitude for understanding animal sounds or complex music (Galantucci et al., 2006).

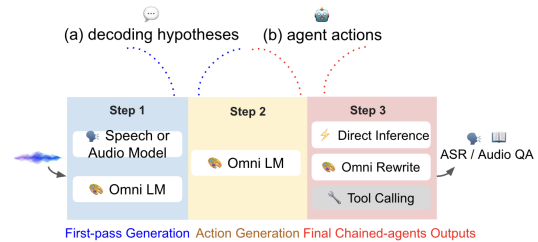


Figure 1: *Speech-Hands* acts as a dynamic orchestrator that predicts a special action token to govern its cognitive strategy.

Inspired by developmental psychology (Selman and Byrne, 1974), we draw a parallel to the human capacity for self-reflection, where children mature from a purely egocentric viewpoint to a stage of *self-reflective perspective-taking*, which serves the critical ability to “step outside” one’s own thoughts, evaluate one’s beliefs against others’, and importantly, to recognize the boundaries of one’s own knowledge. In contrast, current models often operate egocentrically, implicitly trusting their internal perception without the capacity to critically assess its reliability or seek external assistance when necessary. We aim to instill a form of computational self-reflection (Nelson, 1990) into an omni-modal agent, designing a collaborative framework that explicitly reasons about when to trust its own perception, when to defer to an expert, and even when to utilize tools.

We frame these voice understanding models not just as passive predictors, but as agentic that have access to multiple internal and external information sources, and must decide how to best use them. For such an agent, a central decision-making challenge arises (Lebiere and Anderson, 2011): should it rely on its own auditory perception, or consult external suggestions, such as ASR alternatives or other perceptual consultants? Prior work, such as ASR and large language model (LLM) cascaded approach of Generative Error Correction (GER) (Yang et al.,

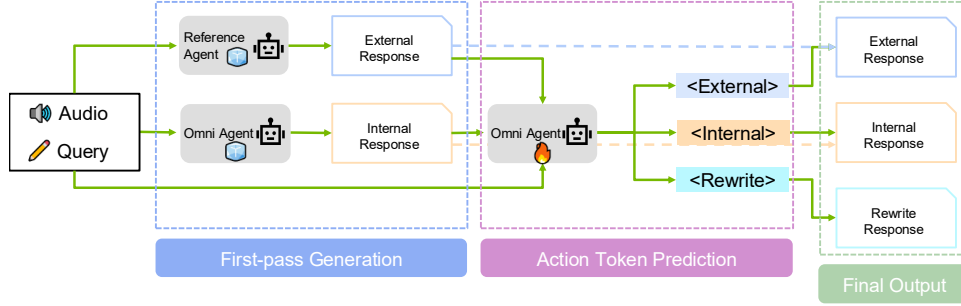


Figure 2: Overview of our proposed Self-Reflection Multimodal GER framework. A special token is generated at the beginning to decide whether to use audio perception (i.e., transcription hypotheses or caption) or not.

2023a), as shown in Figure 5, sidesteps this agentic dilemma entirely. By operating only on text hypotheses without access to the original audio, these methods are fundamentally non-agentic; they cannot weigh internal perception against external advice because they have no internal perception to begin with. Our preliminary experiments reveal that naively combining modalities often degrades performance, as the model struggles to resolve conflicts between its own perception and flawed external suggestions (Kaiser et al., 2021). Without a mechanism to decide which source to trust, the model is easily confused.

To address this, we introduce *Speech-Hands*, a learnable framework that instantiates self-reflection as a core control primitive. As illustrated in Figure 1, our agent operates as a dynamic orchestrator. It begins by aggregating multi-source decoding hypotheses (a) from the first-pass generation. Instead of blindly fusing inputs, the agent critically evaluates them to predict an explicit agent action (b) during the action generation phase. This control token effectively dictates the model’s cognitive strategy: triggering fast direct inference when confidence is high (selecting whether its internal perception or external perceptions), engaging in omni rewrite over available evidence, or initiating tool calling when special utilities are required. In this work, we will focus on the actions of direct inference and omni rewrite, leaving the tool calling action in future work. This approach unifies transcription and reasoning under a single, controllable framework that knows when to trust, when to rethink, and when to ask for help.

### 1.1 Preliminary: The Surprising Failure of Multimodal Correction with Omni-LM

A natural hypothesis is that providing an omni model with both audio and text hypotheses during

supervised fine-tuning (SFT) should enhance GER performance. We tested this assumption by fine-tuning Qwen2.5-Omni (Xu et al., 2025b) to correct N-best hypotheses ( $N = 5$ ) from Whisper-v2-large (Radford et al., 2023) on OpenASR datasets.

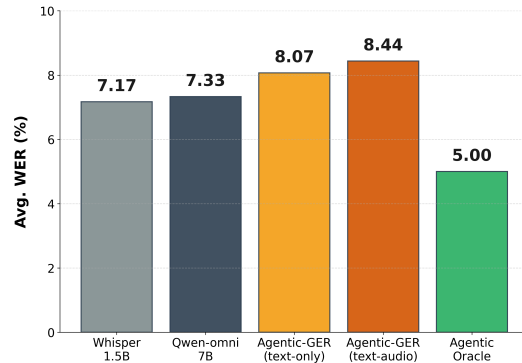


Figure 3: Preliminary results on the cascaded agentic Qwen-omni baseline for generative error correction (GER) with supervised fine-tuning show that both text-only and text-audio GER degrade ASR performance, where the best ASR and LLM combination achieves a low agentic output oracle of 5% WER.

Figure 3 results, however, are surprisingly negative. As shown in our preliminary study (Table 1), this naive SFT approach consistently degrades performance across seven ASR benchmarks, yielding a higher Word Error Rate (WER) than either of the baseline models alone.

To verify that this degradation was not due to suboptimal prompting, we extensively tested varying instructions during SFT, which aims to emphasize internal audio perception, external transcripts, or a balanced fusion. However, as detailed in Appendix A.3, all prompting setups failed to recover performance (e.g., WER increased to 8.52%–9.05% compared to baselines). Furthermore, our zero-shot analysis reveals that the base model lacks intrinsic arbitration capabilities: its decisions are

highly sensitive to prompt wording rather than ground truth, often collapsing into trivial heuristics (see Appendix A.4).

These findings demonstrate a fundamental flaw in the naive omni-LM fusion approach: without a mechanism to adjudicate between its own perception and potentially flawed external advice, the omni-model is easily confused, often amplifying hallucinations or overcorrections as shown in the ASR failure case (Appendix A and Section 6.3). This provides strong motivation for a more principled mechanism that allows the model to learn when and how to incorporate external information.

## 2 Related Work

### 2.1 Voice Retrieval and Agentic Framework

Voice retrieval has become a useful component in augmenting audio tasks such as captioning (Koizumi et al., 2020; Zhao et al., 2023), audio-to-text generation (Huang et al., 2023), dialogue system (Chen et al., 2025), and music generation (Gonzales and Rudzicz, 2024). Recently, voice retrieval modules have been incorporated into multimodal agents (Yang et al., 2023b; Zhang et al., 2024; Wang et al., 2025), providing access to external memory or specialized tools across modalities.

Yet, even in these agentic frameworks, retrieval is often treated as an auxiliary enhancer, rather than as a distinct source of information. When the retrieved content diverges from the internal prediction of the model, current systems lack principled mechanisms for arbitration.

### 2.2 Omni-Modality and Self-Reflection in Multimodal Models

By weaving together text, vision, and audio into a single fabric of understanding, these systems begin to approach the fluidity of human perception (Xu et al., 2025b,c; Goel et al., 2024; Abouelenin et al., 2025). Modality-specific reflection methods (Hu et al., 2025) suggest that introspection within each sensory channel can partially bridge these gaps, aligning representations with quiet precision.

Recent works also introduce explicit self-reflection (Renze and Guven, 2024; Madaan et al., 2023) into multimodal reasoning (Cheng et al., 2024; Fang et al., 2025). The self-reflected video reasoner (Song et al.) iteratively critiques its own visual understanding to reinforce policy stability, while other efforts call for reflective checks against overconfidence and modality neglect (Yang et al.,

2025b). Still, these mechanisms operate *after* perception, which treats reflection as a corrective mirror once fusion has already occurred.

Rather than reflecting on the output, our Speech-Hands framework reflects on the act of perception itself. It learns an action mechanism that decides, in real time, whether to trust its own “ears” or the “words” of others. We aim to discover self-reflection from a post-hoc repair strategy into a preemptive act of perceptual discernment toward an early glimmer of meta-cognition within multimodal understanding.

## 3 Methodology

We present **Speech-Hands**, a learnable omni-agentic framework for audio understanding and reasoning (Figure 2). It enables a multimodal language model to explicitly choose a special token: trusting its own internal perception or defer to external hypotheses, enables efficient **Fast Inference**, while rewriting a new response engages the **Omni Rewrite** process for deeper reasoning. This token is generated during inference and guides the downstream generation process, allowing interpretable and agentic decision-making.

### 3.1 Task Formulation

We formulate a unified self-reflection agent that generalizes across both speech recognition and audio question answering. Given an input audio  $A$ , an optional query  $Q$ , the agent first generates its own response  $H_{\text{omni}}$ , and then combined with an external response  $H_{\text{ext}}$  provided by an external model.

Next, rather than fusing these sources implicitly as normal GER researches, Speech-Hands introduces a learned policy to explicitly choose among them. Our agent model first emits a special action token from the set  $\{\langle\text{internal}\rangle, \langle\text{external}\rangle, \langle\text{rewrite}\rangle\}$  to indicate whether to trust itself or rely on external hypothesis or even whether rewriting a new response after rethinking the audio task and all input resources for the final answer (GER). This self-reflection decision is made based on the full context  $(A, Q, H_{\text{omni}}, H_{\text{ext}})$ , and the selected action conditions the final generation.

To supervise the self-reflection mechanism, we construct ground-truth action token labels on the whole training dataset by comparing the performance of internal, external, and GER predictions. We leverage detailed strategies for action token

229	construction in ASR and audio QA.	
230	<b>3.2 Action Token Construction for ASR</b>	
231	In the ASR setting, we use WER as a pointer to	
232	decide our actions. For each audio example, we	
233	first prompt the omni model to generate the tran-	
234	script $T_{\text{int}}$ and in parallel leverage an ASR model	
235	to predict the external $T_{\text{ext}}$ , we then let omni model	
236	to generate again based on both the audio and the	
237	external $T_{\text{ext}}$ to acquire the GER prediction $T_{\text{ger}}$ .	
238	Next, we compute the WER between the ground-	
239	truth transcript $T_{\text{gt}}$ and each of the three candidates.	
240	If $T_{\text{int}}$ is the same as $T_{\text{gt}}$ ( $WER = 0$ ) or has the	
241	lowest WER, the label is assigned as <code>&lt;internal&gt;</code> ,	
242	this is to encourage the model to trust itself when	
243	it can successfully solve the problem, otherwise	
244	<code>&lt;external&gt;</code> if $T_{\text{ext}}$ performs best, or <code>&lt;rewrite&gt;</code> if	
245	$T_{\text{ger}}$ performs best.	
246	<b>3.3 Action Token Construction for Audio QA</b>	
247	Unlike the ASR setting where token selection is	
248	based on fine-grained WER scores, Audio QA	
249	presents a discrete supervision signal: each predic-	
250	tion is either correct or incorrect. For each instance	
251	consisting of an audio segment $A$ , a question $Q$ ,	
252	and answer choices $\mathcal{C}$ , we first prompt the omni	
253	model to produce an internal prediction $c_{\text{int}}$ based	
254	on $(A, Q)$ . In parallel, we obtain an external pre-	
255	diction $c_{\text{ext}}$ from an audio reasoning model.	
256	We then compare both predictions against the	
257	ground-truth answer $c^*$ . If $c_{\text{int}} = c^*$ , we assign the	
258	label <code>&lt;internal&gt;</code> , encouraging self-reliance when	
259	the model performs well. If $c_{\text{int}}$ fails but $c_{\text{ext}} = c^*$ ,	
260	we assign <code>&lt;external&gt;</code> to delegate control. Other-	
261	wise, when both predictions are incorrect, the label	
262	is set to <code>&lt;rewrite&gt;</code> , signaling a need to re-evaluate	
263	the question with all available context.	
264	However, this binary decision process introduces	
265	instability during training. External predictions	
266	can be stochastic. Therefore, repeated sampling	
267	may yield different answers, especially under high	
268	uncertainty or directly change another external	
269	model can also lead to a different accuracy. This	
270	stochasticity makes the decision boundary between	
271	<code>&lt;external&gt;</code> and <code>&lt;rewrite&gt;</code> inherently less robust.	
272	To mitigate this, we adopt a multiple decoding-	
273	based strategy: for each example, we sample the	
274	external model five times and collect their predicted	
275	choices. If the majority of predictions match the	
276	ground-truth answer, we assign <code>&lt;external&gt;</code> ; oth-	
277	erwise, we assign <code>&lt;rewrite&gt;</code> . This approach stabi-	
	lizes supervision by reducing the variance in exter-	278
	nal outputs and yields more reliable action labels.	279
	<b>3.4 Prompt Formatting and Training</b>	280
	Subsequently, each training instance is formatted	281
	as a single target string consisting of the decision	282
	token followed by the final target transcript or an-	283
	swer, e.g. <code>&lt;rewrite&gt;</code> + ground-truth transcription.	284
	This unified string allows the model to learn not	285
	only how to generate the task output but also how	286
	to decide the action before generation. We adopt an	287
	instruction-style prompt to guide the model to make	288
	decisions as shown in Appendix B. During train-	289
	ing, we optimize a single cross-entropy loss over	290
	the concatenated target sequence, which jointly	291
	supervises the action token and the subsequent pre-	292
	diction. Concretely, the model first predicts the	293
	action token and then continues decoding the target	294
	transcript or answer; both parts contribute to the	295
	same loss. This end-to-end objective encourages	296
	the model to internalize the mapping from multi-	297
	modal evidence to action choice and to generate	298
	the corresponding output under that decision.	299
	<b>3.5 Agentic Inference via Action Tokens</b>	300
	At inference time, the model receives the same	301
	multimodal inputs as during training. The model	302
	performs a two-stage decoding process: it first	303
	emits an action token: <code>&lt;internal&gt;</code> , <code>&lt;external&gt;</code> ,	304
	or <code>&lt;rewrite&gt;</code> . The decoding process is subject to	305
	decide which information source to prioritize, and	306
	then generates the final output accordingly. This	307
	explicit agentic self-reflection mechanism offers	308
	interpretability and control over how the model bal-	309
	ances internal perception and external knowledge.	310
	It enables direct analysis (e.g., F1 score) of when	311
	the model relies on its own understanding, consul-	312
	itants from external systems, or synthesizes a new	313
	response. The unified prediction format ensures	314
	that the model not only learns what to generate, but	315
	also which to trust across both speech recognition	316
	and audio reasoning tasks.	317
	<b>4 Experimental Setup</b>	318
	<b>4.1 Datasets</b>	319
	<b>Speech Recognition.</b> We use seven represen-	320
	tative datasets covering a range of domains,	321
	styles, and noise conditions from OpenASR leader-	322
	board: AMI (Carletta, 2007) (meeting speech),	323
	Tedlium (Hernandez et al., 2018) (TED talks), Gi-	324
	gaSpeech (Chen et al., 2021) (large-scale podcasts	325

Dataset	AMI	Tedlium	Gigspeech	Spgispeech	VoxPopuli	Libri-clean	Libri-other	avg. WER ↓
<b>ASR model or Omni-LLM</b>								
Whisper-v2-large	16.88	4.32	11.45	3.94	7.57	2.91	5.15	7.17
Canary-1b-v2	19.80	4.78	11.66	3.08	6.35	1.73	3.17	7.22
Parakeet-tdt-0.6b-v3	12.69	4.90	12.24	3.16	6.48	1.89	3.37	6.68
Qwen2.5_omni	19.77	5.17	11.26	4.58	6.59	2.09	3.85	7.33
Phi-4-MM	11.69	<b>2.90</b>	<b>9.78</b>	3.13	<b>5.93</b>	1.68	3.83	6.14
Gemini-2-Flash	21.58	3.01	10.71	3.82	7.89	2.49	5.84	8.56
GPT-4o-voice	57.76	5.79	13.64	5.66	10.83	3.48	7.97	15.76
<b>Qwen2.5-Omni: Cascaded</b>								
GER: ⇒ Whisper-v2-large	23.44	6.15	12.15	3.94	7.53	2.97	4.89	8.44
GER: ⇒ canary	24.58	6.38	12.43	4.02	7.72	3.05	5.01	8.74
GER: ⇒ parakeet	22.91	6.09	12.10	3.98	7.49	2.92	4.84	8.33
<b>Qwen2.5-Omni: Parallel</b>								
Speech-Hands ⇌ whisper-v2	15.03	4.45	12.37	3.01	6.49	1.86	3.46	6.67
Speech-Hands ⇌ canary	15.29	4.21	10.87	<b>2.17</b>	5.96	<b>1.61</b>	<b>3.07</b>	6.17
Speech-Hands ⇌ parakeet	<b>11.20</b>	4.37	11.10	2.26	6.02	1.67	3.18	<b>5.69</b>

Table 1: WER (%) results across 7 datasets, with the average WER shown in the rightmost column. Speech-Hands training significantly outperforms both baseline systems (ASR model, Qwen) and prior cascaded GER setups.

and YouTube-style speech), SPGISpeech (O’Neill et al., 2021) (long-form read speech), VoxPopuli (Wang et al., 2021) (English subset of multilingual political recordings), and LibriSpeech (Panayotov et al., 2015) (clean and noisy audiobook speech). For baseline fine-tuning and prompt-based GER, we use all available training sets. For our proposed method, unless otherwise specified (w/ FULL Datasets), we train on at most 20,000 examples per dataset, this is due to the limitation of heavy computation requirements when doing inference on the whole training set for internal, external and GER (token distribution is discussed in 6.1).

**Audio Reasoning.** We evaluate on the multi-domain audio question-answering benchmark (Yang et al., 2025a) (MD-Audio), which consists of multiple-choice questions grounded in real audio clips. This benchmark includes three complementary subsets that probe different reasoning capabilities. While previous audio reasoning benchmarks such as MMAU (Sakshi et al., 2024a) and MMAR (Ma et al., 2025) provide only MMLU-style test sets, MD-Audio releases both training and development sets, enabling evaluation of the proposed trainable agentic framework, as shown in Table 8. Detailed dataset descriptions can be found in Appendix C. For each sample, we construct inputs as described in Section 3, including the original audio, a first-pass internal prediction from Omni model, one external hypothesis and GER result (for ASR).

## 5 Results

### 5.1 Training Details

All experiments are conducted using the Qwen2.5-Omni model. We extend its tokenizer to include three special action tokens, <internal>, <external> and <rewrite>, used during both training and inference. Models are trained using the standard supervised fine-tuning (SFT) objective with a cross-entropy loss. We train for 5 epochs with fp16. The batch size is set to 64, and the learning rate is initialized at 1e-4 with cosine decay. All experiments adopt greedy decoding.

### 5.2 Baselines

**External ASR models:** We include three high-performing supervised speech recognition models as external references: Whisper-v2-large (Radford et al., 2022), Canary-1B-v2, and Parakeet-TDT-0.6B-v3 (Sekoyan et al., 2025). These systems operate in a closed transcription setting and provide non-generative references.

**External audio QA model:** For multi-choice audio QA, we include Audio Flamingo 3 (Goel et al., 2025), a speech-language model with audio frontend capabilities serving as a strong baseline.

We also include the latest model in comparison, e.g., Phi-4-MM (Microsoft et al., 2025), Gemini-2-Flash, GPT-4o-voice in ASR evaluations.

### 5.3 ASR Results

Table 1 presents WER results across seven diverse datasets, and we find that: (1) Speech-Hands outperforms all baselines. Furthermore, our approach with strong ASR models (canary and parakeet)

Model / Setting	Bio-acoustic	Soundscape	Complex QA	avg. Acc. ↑
<b>Audio LM or Omni Model</b>				
Gemini-2-Flash	42.03	46.34	59.89	56.61
Qwen2.5-Omni	47.32	56.32	59.89	57.87
AudioFlamingo 3 (AF3)	71.88	57.31	81.26	74.49
<b>Qwen2.5-Omni Baselines</b>				
+ SFT with official training data	78.13	34.65	76.61	63.13
+ GRPO with official training data	78.09	39.43	79.12	65.54
+ GRPO with external audio data (Li et al., 2025a)	62.32	<b>72.10</b>	82.15	75.10
GER: ⇒ AF3 (cascaded agentic)	76.29	52.02	77.48	68.93
<b>Qwen2.5-Omni: Speech-Hands</b>				
⇒ (parallel agentic): SFT with official training data	67.86	58.29	83.34	75.75
⇒ (parallel agentic) + majority sampling	<b>81.25</b>	59.4	<b>85.7</b>	<b>77.37</b>

Table 2: AudioQA and acoustic content reasoning accuracy (%) across knowledge-intensive bioacoustic QA (Sayigh et al., 2016), multi-sound-object soundscapes, and MMAU-style (Sakshi et al., 2024a) complex audio QA tasks.

389 achieves the lowest WER, even with only 20k training  
390 examples, outperforming both ASR models  
391 or Omni-LLMs; (2) The prompt GER over Whisper  
392 lags significantly behind token-based methods.  
393 This underscores the importance of explicit control  
394 via action tokens rather than relying solely on  
395 natural-language prompts; (3) While pre-trained  
396 models like Whisper and Qwen perform well on  
397 curated datasets such as LibriSpeech-clean, their  
398 performance degrades significantly on conversational  
399 benchmarks like AMI. Notably, despite Qwen’s  
400 relatively weaker base ASR performance, our  
401 framework enhances its generalization to the extent  
402 that it surpasses stronger baselines such as Phi-4-MM  
403 on the average performance, demonstrating the  
404 stability and transferability of Speech-Hands on  
405 both clean and noisy datasets.

## 406 5.4 AudioQA Results

407 Table 2 reports accuracy on each sub-task and  
408 overall average accuracy: (1) Our final setup  
409 (Speech-Hands + majority sampling) achieves the  
410 highest average accuracy (77.37%), outperforming  
411 all baselines and pre-trained models. It performs  
412 particularly well on Complex QA (85.70%) and  
413 Bioacoustics QA (81.25%), indicating its ability  
414 to handle both abstract reasoning and fine-grained  
415 audio patterns; (2) Standard supervised fine-tuning  
416 (SFT) and prompt-based GER exhibit mixed  
417 results. SFT achieves good accuracy on Bioacoustics  
418 (78.13%) but fails on Soundscapes (34.65%).  
419 Prompt-based GER also fails on both Soundscape  
420 and Complex QA compared with flamingo 3  
421 baseline. These results highlight the robustness  
422 of our agentic framework in diverse audio reasoning  
settings.

## 423 6 Additional Analysis

### 424 6.1 Accuracy of Action Token Prediction

425 We analyze the model’s ability to correctly emit the  
426 three action tokens of <internal>, <external>, and  
427 <rewrite>, in the ASR setting under the Action  
428 Tokens + Whisper configuration. These tokens  
429 interpret whether the model is making correct  
430 decisions to guide its final prediction.

431 Table 3 shows both the training distribution and  
432 the test-time precision, recall, and F1 scores. The  
433 distribution highlights a strong internal bias: across  
434 all datasets, <internal> dominates, exceeding  
435 95% in Libri-clean, Libri-other, spgispeech, and  
436 Voxpopuli. In contrast, <external> is sparsely  
437 supervised, where below 1% in Librispeech and  
438 <rewrite> is extremely rare everywhere, often less  
439 than 2%. This imbalance poses a natural challenge  
440 for learning reliable action.

441 Despite this skew, the model demonstrates robust  
442 performance for the two tokens. On test data,  
443 <internal> predictions achieve F1 scores above  
444 0.8 on most datasets (0.91 on spgispeech, 0.94  
445 on Libri-clean, 0.90 on Libri-other), indicating  
446 that the model can reliably recognize when its  
447 own decoding is sufficient. Even for the much  
448 rarer <external> token, the model attains high  
449 F1 scores, showing strong generalization of  
450 deferring to external hypotheses despite limited  
451 supervision.

452 The <rewrite> token proves to be the most  
453 challenging, with F1 scores below 0.4 in all but  
454 one dataset and zero in Librispeech, where  
455 positive training examples are extremely rare.  
456 A closer examination reveals that precision  
457 consistently exceeds recall, indicating that when  
458 the model does emit <rewrite>, its decision is  
generally correct but under-triggered in the omni  
model. This sug-

Dataset	AMI	Tedlium	Gigaspeech	SPGIspeech	Voxpopuli	Libri-clean	Libri-other
<b>Training Distribution</b>							
<internal>	67.95%	86.48%	87.76%	96.25%	93.73%	98.96%	98.96%
<external>	31.01%	11.18%	11.8%	3.64%	6.09%	0.96%	0.96%
<rewrite>	1.04%	2.34%	0.44%	1.21%	0.18%	0.1%	0.1%
<b>Test Distribution</b>							
<internal>	70.28%	83.57%	85.94%	95.42%	92.68%	98.92%	98.75%
<external>	26.91%	15.12%	12.41%	3.81%	6.47%	0.98%	1.08%
<rewrite>	2.27%	1.31%	1.65%	0.77%	0.85%	0.1%	0.17%
<b>&lt;internal&gt; on Test</b>							
Precision	0.85	0.63	0.77	0.89	0.85	0.88	0.83
Recall	0.78	0.72	0.90	0.94	0.90	0.99	0.99
F1	0.81	0.67	0.83	0.91	0.87	0.94	0.9
<b>&lt;external&gt; on Test</b>							
Precision	0.88	0.96	0.83	0.82	0.71	0.81	0.72
Recall	0.89	0.81	0.78	0.76	0.60	0.73	0.75
F1	0.89	0.88	0.80	0.79	0.65	0.77	0.74
<b>&lt;rewrite&gt; on Test</b>							
Precision	0.62	0.33	0.32	0.52	0.50	0.0	0.0
Recall	0.24	0.21	0.05	0.36	0.21	0.0	0.0
F1	0.39	0.28	0.08	0.43	0.33	0.0	0.0

Table 3: Training distribution and test-time F1 scores for Speech-Hands’ action tokens.

gests a cautious yet reasonably reliable rewrite detector, whose coverage could be further improved through targeted data augmentation.

Overall, these results validate the effectiveness of the proposed agentic action: even under heavy class imbalance, the model learns to accurately identify when to trust its own predictions versus when to consult external information. The main bottleneck remains the <rewrite> case, suggesting that richer sampling or augmentation strategies may be needed to stabilize this decision in future work.

## 6.2 Confusion Analysis In AudioQA

Subset	<in>/<ex>/<re>
Bio-acoustic QA	106/75/43
Soundscape QA	343/185/81
Complex QA	978/555/100

Table 4: Oracle Statistics of the three DCASE2025 AudioQA test sets used in our experiments.

We analyze the confusion matrix over the three subsets under the w/ multiple sampling setup (the oracle token distribution is shown in Table 4). As shown in Figure 4, the confusion between <internal> and <external> remains relatively low, suggesting that the model can effectively distinguish between them, especially in the Complex QA subset. However, the Soundscape subset shows slightly increased overlap, possibly due to the difficulty of Soundscape compared with other parts, also their original performances (Qwen 56.32 v.s.

Flamingo 57.31) are quite closed, leading to much smaller sets of <external>. Besides, the highest confusion still lies with the <rewrite> class. As summarized in Table 8, the number of training tokens labeled as <rewrite> is significantly smaller across all subsets. Such sparsity likely limits the model’s tendency to generalize the rewriting behavior during generation, but when it generates <rewrite>, the accuracy is quite high and robust as shown in the Figure. These findings highlight that even with imbalanced actions in training, the current F1 scores can still reliably improve the final performance, emphasizing the effectiveness of action tokens.

### AudioQA Inference-Time Case Study

**Q.** Based on the audio, which natural phenomenon could be occurring?

- A. Earthquake B. Thunderstorm  
C. Forest fire D. Snowstorm

**Internal pred:** B. Thunderstorm

**External pred:** B. Thunderstorm

**Pred:** <rewrite> C. Forest fire (✓)

### ASR Inference-Time Case Study

**Internal:** you in the way marguerite but how

**External:** you ll in the way marguerite but how

**Rewrite:** you are in the way marguerite but how

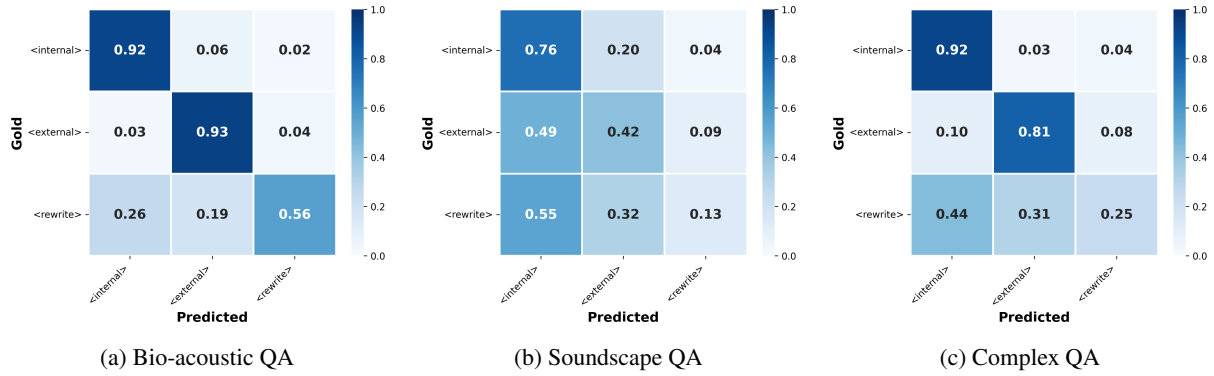


Figure 4: Confusion matrices of Speech-Hands’ agentic action execution for audio QA and reasoning three subsets on (a) bio-acoustic QA, (b) temporal and sound event QA, and (c) complex audio information QA.

**Pred:** <internal> you in the way mar-  
guerite but how (✓)

### 6.3 Category Analysis and Case Study

To provide systematic insight beyond individual cases, during our preliminary, we categorize three primary failure modes observed in ASR when Speech-Hands is absent. These categories further motivate our design choice of action tokens: (1) **External-induced Misguidance:** the internal omni-model is correct, but the model is misled by an error-heavy external hypothesis. This is the most damaging failure mode and directly motivates the action-token arbitration design; (2) **Overcorrection:** LLMs tend to “complete” disfluent or partially heard speech, introducing hallucinated insertions or semantic expansions; (3) **Undercorrection in Dual-Failure:** when both internal and external predictions are wrong, the omni model often selects one erroneous hypothesis instead of generating a rewrite, due to the lack of an explicit rewrite mechanism. These observations align with the high-precision behavior of the <rewrite> token and explain the improvements: the model selectively activates <rewrite> only when necessary.

We present illustrative examples (more in Appendix E) to show how the model makes selective use of internal, external, and rewrite hypotheses across modalities. In the AUDIOQA CASE), both the internal and external models predict “B. Thunderstorm”, likely influenced by surface acoustic features such as low-frequency rumbling. The rewrite path, however, generates “C. Forest fire”, which aligns with the ground truth, demonstrating the strong influence of <rewrite> that even two

models have the same first-pass prediction.

In contrast, the ASR CASE reveals a different decision dynamic. Although the rewrite produces a more fluent variant, the model opts to retain the baseline hypothesis, judging the original phrase as sufficiently accurate. This indicates the model’s ability to avoid overcorrection issue in prior GER researches when input ambiguity is low. These examples underscore the benefits of explicit action tokens: the model can either rely on internal or external model predictions or revise them when necessary, yielding both flexibility and robustness in audio tasks.

## 7 Conclusion

In this work, we proposed a learnable voice-agentic framework Speech-Hands for teaching omni models when to trust itself versus when to consult external audio perception. By casting the problem with explicit <internal>, <external>, and <rewrite> action tokens, our experimental results across AudioQA and ASR benchmarks demonstrate strong performance improvements beyond strong baselines, especially when direct finetuning and GER training fail, Speech-Hands can still robustly generate the best prediction.

This framework also benefits the interpretability in analysis, the model achieves high F1 scores for both <internal> and <external> tokens, even under imbalanced training conditions. While the <rewrite> token is rarer, its precision notably exceeds recall, indicating that the model can accurately identify necessary rewrites when it does trigger them. Overall, our method offers an effective framework to inject explicit actions into agent decision, toward reliable audio intelligence.

## 566 Limitations

567 Despite promising results, our study presents sev-  
568 eral limitations that offer avenues for future explo-  
569 ration.

570 **Token imbalance and rewrite sparsity.** Our  
571 training data exhibits an inherent imbalance  
572 across action tokens (<internal>, <external>,  
573 <rewrite>). While both <internal> and  
574 <external> achieve high F1 scores, <rewrite>  
575 remains under-trained on many datasets. This spar-  
576 sity partly reflects that certain audio QA datasets  
577 rarely require rewriting but this contextual infor-  
578 mation sparsity also reveals a modeling challenge.  
579 Future work may explore more principled strate-  
580 gies for balancing token distribution or adaptively  
581 reshaping decision boundaries, especially under  
582 varying persona settings or task configurations.

583 **Limited ASR training subset.** Our current ASR  
584 experiments are trained on a restricted subset of  
585 data. While the model already achieves strong  
586 performance, it likely underutilizes the available  
587 signal. Scaling up training with larger ASR datasets  
588 or augmenting with synthetic audio variants may  
589 unlock further gains.

590 **No exploration of transfer or multi-external se-**  
591 **tups.** We do not yet study transfer capabilities.  
592 For example, training with one external ASR model  
593 and testing with another. Moreover, our current  
594 system only accepts a single external prediction.  
595 Extending the framework to handle multiple ex-  
596 ternal models, each represented by distinct deci-  
597 sion tokens, could significantly improve robustness  
598 and enable broader deployment across diverse real-  
599 world pipeline

## 600 References

601 Abdelrahman Abouelenin, Atabak Ashfaq, Adam Atkin-  
602 son, Hany Awadalla, Nguyen Bach, Jianmin Bao,  
603 Alon Benhaim, Martin Cai, Vishrav Chaudhary,  
604 Congcong Chen, et al. 2025. Phi-4-mini techni-  
605 cal report: Compact yet powerful multimodal lan-  
606 guage models via mixture-of-loras. *arXiv preprint*  
607 *arXiv:2503.01743*.

608 Sharath Adavanne, Archontis Politis, and Tuomas Vir-  
609 tanen. 2019. A multi-room reverberant dataset for  
610 sound event localization and detection. *Preprint*,  
611 arXiv:1905.08546.

612 A. Calcus. 2024. Development of auditory scene analy-  
613 sis: a mini-review. *Frontiers in Human Neuroscience*,  
614 18:1352247.

Jean Carletta. 2007. Unleashing the killer corpus: ex-  
periences in creating the multi-everything ami meet-  
ing corpus. *Language Resources and Evaluation*,  
41(2):181–190.

Guoguo Chen, Shuzhou Chai, Guan-Bo Wang, Jiayu  
Du, Wei-Qiang Zhang, Chao Weng, Dan Su, Daniel  
Povey, Jan Trmal, Junbo Zhang, Mingjie Jin, Sanjeev  
Khudanpur, Shinji Watanabe, Shuaijiang Zhao, Wei  
Zou, Xiangang Li, Xuchen Yao, Yongqing Wang,  
Zhao You, and Zhiyong Yan. 2021. *Gigaspeech:*  
*An evolving, multi-domain asr corpus with 10,000*  
*hours of transcribed audio*. In *Interspeech 2021*, in-  
terspeech\_2021. ISCA.

Yifu Chen, Shengpeng Ji, Haoxiao Wang, Ziqing Wang,  
Siyu Chen, Jinzheng He, Jin Xu, and Zhou Zhao.  
2025. *Wavrag: Audio-integrated retrieval augmented*  
*generation for spoken dialogue models*. *Preprint*,  
arXiv:2502.14727.

Kanzhi Cheng, Yantao Li, Fangzhi Xu, Jianbing Zhang,  
Hao Zhou, and Yang Liu. 2024. Vision-language  
models can self-improve reasoning via reflection.  
*arXiv preprint arXiv:2411.00855*.

Yangui Fang, Baixu Cheng, Jing Peng, Xu Li, Yu Xi,  
Chengwei Zhang, and Guohui Zhong. 2025. Fewer  
hallucinations, more verification: A three-stage llm-  
based framework for asr error correction. *arXiv*  
*preprint arXiv:2505.24347*.

B. Galantucci, C. A. Fowler, and M. T. Turvey. 2006.  
*The motor theory of speech perception reviewed*.  
*Psychonomic Bulletin & Review*, 13(3):361–377. Er-  
ratum in: *Psychon Bull Rev*. 2006 Aug;13(4):742.

Arushi Goel, Sreyan Ghosh, Jaehyeon Kim, Sonal Ku-  
mar, Zhifeng Kong, Sang gil Lee, Chao-Han Huck  
Yang, Ramani Duraiswami, Dinesh Manocha, Rafael  
Valle, and Bryan Catanzaro. 2025. *Audio flamingo 3:*  
*Advancing audio intelligence with fully open large*  
*audio language models*. *Preprint*, arXiv:2507.08128.

Arushi Goel, Karan Sapra, Matthieu Le, Rafael Valle,  
Andrew Tao, and Bryan Catanzaro. 2024. Omcat:  
Omni context aware transformer. *arXiv preprint*  
*arXiv:2410.12109*.

Robie Gonzales and Frank Rudzicz. 2024. A retrieval  
augmented approach for text-to-music generation. In  
*Proceedings of the 3rd Workshop on NLP for Mu-*  
*sic and Audio (NLP4MusA)*, pages 31–36, Oakland,  
USA. Association for Computational Linguistics.

François Hernandez, Vincent Nguyen, Sahar Ghan-  
nay, Natalia Tomashenko, and Yannick Estève. 2018.  
*TED-LIUM 3: Twice as Much Data and Corpus*  
*Repertition for Experiments on Speaker Adaptation*,  
page 198–208. Springer International Publishing.

Rui Hu, Delai Qiu, Shuyu Wei, Jiaming Zhang, Yining  
Wang, Shengping Liu, and Jitao Sang. 2025. In-  
vestigating and enhancing vision-audio capability in  
omnimodal large language models. In *Findings of*  
*the Association for Computational Linguistics: ACL*



789	Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeef Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pocrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. <a href="#">Gpt-4 technical report</a> . <i>Preprint</i> , arXiv:2303.08774.	
834	Patrick K. O’Neill, Vitaly Lavrukhin, Somshubra Majumdar, Vahid Noroozi, Yuekai Zhang, Oleksii Kuchaiev, Jagadeesh Balam, Yuliya Dovzhenko, Keenan Freyberg, Michael D. Shulman, Boris Ginsburg, Shinji Watanabe, and Georg Kucsko. 2021. <a href="#">Spgispeech: 5,000 hours of transcribed financial audio for fully formatted end-to-end speech recognition</a> . In <i>Interspeech 2021</i> , pages 1434–1438.	
842	Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. <a href="#">Librispeech: An asr corpus based on public domain audio books</a> . In <i>2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 5206–5210.	
847	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. <a href="#">Robust speech recognition via large-scale weak supervision</a> . <i>Preprint</i> , arXiv:2212.04356.	
851	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. <a href="#">Robust speech recognition via large-scale weak supervision</a> . In <i>International conference on machine learning</i> , pages 28492–28518. PMLR.	851 852 853 854 855
	Matthew Renze and Erhan Guven. 2024. <a href="#">Self-reflection in llm agents: Effects on problem-solving performance</a> . <i>arXiv preprint arXiv:2405.06682</i> .	856 857 858
	S Sakshi, Utkarsh Tyagi, Sonal Kumar, Ashish Seth, Ramaneswaran Selvakumar, Oriol Nieto, Ramani Duraiswami, Sreyan Ghosh, and Dinesh Manocha. 2024a. <a href="#">Mmau: A massive multi-task audio understanding and reasoning benchmark</a> . <i>arXiv preprint arXiv:2410.19168</i> .	859 860 861 862 863 864
	S Sakshi, Utkarsh Tyagi, Sonal Kumar, Ashish Seth, Ramaneswaran Selvakumar, Oriol Nieto, Ramani Duraiswami, Sreyan Ghosh, and Dinesh Manocha. 2024b. <a href="#">Mmau: A massive multi-task audio understanding and reasoning benchmark</a> . <i>Preprint</i> , arXiv:2410.19168.	865 866 867 868 869 870
	Laela Sayigh, Mary Ann Daher, Julie Allen, Helen Gordon, Katherine Joyce, Claire Stuhlmann, and Peter Tyack. 2016. <a href="#">The watkins marine mammal sound database: an online, freely accessible resource</a> . In <i>Proceedings of Meetings on Acoustics</i> , volume 27, page 040013. Acoustical Society of America.	871 872 873 874 875 876
	Monica Sekoyan, Nithin Rao Koluguri, Nune Tadevosyan, Piotr Zelasko, Travis Bartley, Nikolay Karpov, Jagadeesh Balam, and Boris Ginsburg. 2025. <a href="#">Canary-1b-v2 &amp; parakeet-tdt-0.6b-v3: Efficient and high-performance models for multilingual asr and ast</a> . <i>Preprint</i> , arXiv:2509.14128.	877 878 879 880 881 882
	Robert L Selman and Diane F Byrne. 1974. <a href="#">A structural-developmental analysis of levels of role taking in middle childhood</a> . <i>Child development</i> , pages 803–806.	883 884 885 886
	Zihan Song, Xin Wang, Zi Qian, Hong Chen, Longtao Huang, Hui Xue, and Wenwu Zhu. <a href="#">Modularized self-reflected video reasoner for multimodal llm with application to video question answering</a> . In <i>Forty-second International Conference on Machine Learning</i> .	887 888 889 890 891 892
	Changhan Wang, Morgane Riviere, Ann Lee, Anne Wu, Chaitanya Talnikar, Daniel Haziza, Mary Williamson, Juan Pino, and Emmanuel Dupoux. 2021. <a href="#">VoxPopuli: A large-scale multilingual speech corpus for representation learning, semi-supervised learning and interpretation</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 993–1003, Online. Association for Computational Linguistics.	893 894 895 896 897 898 899 900 901 902 903
	Zixuan Wang, Chi-Keung Tang, and Yu-Wing Tai. 2025. <a href="#">Audio-agent: Leveraging llms for audio generation, editing and composition</a> . <i>Preprint</i> , arXiv:2410.03335.	904 905 906 907

908 Zhifei Xie and Changqiao Wu. 2024. [Mini-omni: Language models can hear, talk while thinking in streaming](#). *Preprint*, arXiv:2408.16725.

909

910

911 Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting

912 He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan,

913 Kai Dang, Bin Zhang, Xiong Wang, Yunfei Chu, and

914 Junyang Lin. 2025a. [Qwen2.5-omni technical report](#).  
915 *Preprint*, arXiv:2503.20215.

916 Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting

917 He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan,

918 Kai Dang, et al. 2025b. [Qwen2.5-omni technical](#)  
919 [report](#). *arXiv preprint arXiv:2503.20215*.

920 Jin Xu, Zhifang Guo, Hangrui Hu, Yunfei Chu, Xiong

921 Wang, Jinzheng He, Yuxuan Wang, Xian Shi, Ting

922 He, Xinfa Zhu, et al. 2025c. [Qwen3-omni technical](#)  
923 [report](#). *arXiv preprint arXiv:2509.17765*.

924 Chao-Han Huck Yang, Sreyan Ghosh, Qing Wang,

925 Jaeyeon Kim, Hengyi Hong, Sonal Kumar,

926 Guirui Zhong, Zhifeng Kong, S Sakshi, Vaibhavi

927 Lokegaonkar, Oriol Nieto, Ramani Duraiswami,

928 Dinesh Manocha, Gunhee Kim, Jun Du, Rafael

929 Valle, and Bryan Catanzaro. 2025a. [Multi-domain](#)  
930 [audio question answering toward acoustic content](#)  
931 [reasoning in the dcase 2025 challenge](#). *Preprint*,  
932 arXiv:2505.07365.

933 Chao-Han Huck Yang, Yile Gu, Yi-Chieh Liu, Shalini

934 Ghosh, Ivan Bulyko, and Andreas Stolcke. 2023a.  
935 Generative speech recognition error correction with  
936 large language models and task-activating prompting.  
937 In *2023 IEEE Automatic Speech Recognition and*  
938 *Understanding Workshop (ASRU)*, pages 1–8. IEEE.

939 Qize Yang, Shimin Yao, Weixuan Chen, Shenghao Fu,

940 Detao Bai, Jiaying Zhao, Boyuan Sun, Bowen Yin,

941 Xihan Wei, and Jingren Zhou. 2025b. [Humanomniv2:](#)  
942 [From understanding to omni-modal reasoning with](#)  
943 [context](#). *arXiv preprint arXiv:2506.21277*.

944 Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin

945 Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu,

946 Ce Liu, Michael Zeng, and Lijuan Wang. 2023b.  
947 [Mm-react: Prompting chatgpt for multimodal](#)  
948 [reasoning and action](#). *Preprint*, arXiv:2303.11381.

949 Lu Zhang, Tiancheng Zhao, Heting Ying, Yibo Ma,

950 and Kyusong Lee. 2024. [Omagent: A multi-](#)  
951 [modal agent framework for complex video under-](#)  
952 [standing with task divide-and-conquer](#). *Preprint*,  
953 arXiv:2406.16620.

954 Jinming Zhao, Gholamreza Haffari, and Ehsan Shareghi.

955 2023. [Generating synthetic speech from SpokenVo-](#)  
956 [cab for speech translation](#). In *Findings of the Asso-*  
957 *ciation for Computational Linguistics: EACL 2023*,  
958 pages 1975–1981, Dubrovnik, Croatia. Association  
959 for Computational Linguistics.

## A Failure Case of Multimodal GER

### A.1 Amplifying Hallucination

This is often a situation where external audio perception is not strong enough or N-best decoding introduces noise for GER. Especially when the omni model is encouraged to take all hypotheses in consideration.

#### ASR Failure Case

**Internal:** insane

**External (Whisper 5-best):** [' Gimseeinnnnnn', ' You can say.', ' Insta.', ' Wednesday.', " I'm from Phelps County, I'm gonna see what this guy's doing."]

**GER Pred:** I'm from Phelps County, I'm gonna see what this guy's doing

**Gold:** insane

### A.2 Overcorrection

Due the language modeling nature, LLMs tend to revise the transcription to be more like a complete sentence, which sometimes caused "overcorrection."

#### ASR Failure Case

**Internal:** you in the way marguerite but how

**External (Whisper 5-best):** [' you ll in the way marguerite but how', 'you in the way marguerite but how', 'you in the way marguerite but how.', ' you in the way marguerite but how.', " you in the way marguerite but how."]

**GER Pred:** you are in the way marguerite but how

**Gold:** you in the way marguerite but how

### A.3 Prompt Ablations in Preliminary SFT

We conduct these experiments during early-stage exploration: for inputs we investigated both (a) audio + external whisper 5-best and (b) audio + internal 1-best + external whisper 5-best, together with four prompting strategies instructing the model to emphasize internal hypotheses, external hypotheses, audio grounding, or a balanced fusion. The average WER results on OpenASR are shown in Table 5.

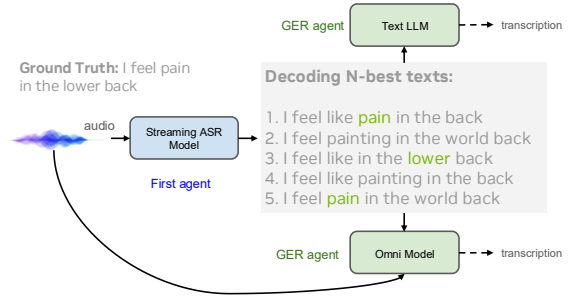


Figure 5: Text-based GER uses only ASR hypotheses. This setup fails to correct deletion or hallucination if all hypotheses are wrong. Multimodal GER include the original audio as grounding to improve error correction.

AVG. WER	(a)	(b)
Emphasize internal	-	8.63
Emphasize external	8.58	8.67
Emphasize audio	9.02	9.05
Balanced	8.44	8.52

Table 5: Prompt ablation results.

### A.4 Zero-shot Qwen in Preliminary

We test three different zero-shot prompting strategies as shown in Table 6. We find that zero-shot decisions are highly prompt-sensitive: the model often collapses into trivial heuristics, while the balanced prompt produces unstable and inconsistent behavior across samples. The corresponding  $2 \times 2$  confusion matrices show large off-diagonal mass for all zero-shot prompts, confirming that zero-shot Qwen does not perform genuine self-reflection.

Table 6: Confusion matrix of zero-shot decisions under different prompting strategies. The results show that the model's arbitration is highly sensitive to prompt wording rather than the ground truth correctness (Oracle), often collapsing into trivial heuristics.

Ground Truth	Internal-biased Prompt		External-biased Prompt		Balanced Prompt	
	Pred. Int	Pred. Ext	Pred. Int	Pred. Ext	Pred. Int	Pred. Ext
Oracle Internal	0.83	0.17	0.35	0.65	0.68	0.32
Oracle External	0.71	0.29	0.14	0.86	0.34	0.66

## B Prompt Template

Below is the prompt template for ASR:

You are an omni-agent for speech understanding with access to three inputs:  
 (1) The original audio;  
 (2) Five transcription hypotheses from another ASR system (external);  
 (3) Your own first-pass transcription (internal).

1005	Your task is to:	about 31 marine mammal species with diverse	1060
1006	- First decide whether your internal	acoustic ranges, habitats, and vocalization dura-	1061
1007	transcription is reliable.	tions. Tasks include species classification, vocal-	1062
1008	- If yes, output <internal> and your	ization type recognition, factual retrieval, interpre-	1063
1009	transcription.	tation of acoustic features, and comparative rea-	1064
1010	- If the external system is more reliable,	soning. The dataset includes approximately 0.7K	1065
1011	output <external> and use one of its	training and 0.2K development QA pairs. Audio	1066
1012	hypotheses.	clips range in sample rate from 600 Hz to 160 kHz	1067
1013	- Otherwise, output <rewrite> and	and in duration from 0.4 s to 625 s, allowing evalu-	1068
1014	generate a new answer using both sources	ation under highly varied acoustic conditions. All	1069
1015	and the audio.	audio is sourced from the Watkins Marine Mam-	1070
1016		mal Sound Database (Woods Hole Oceanographic	1071
1017	Also, here is the prompt for AudioQA, we ex-	Institution; New Bedford Whaling Museum), and	1072
1018	PLICITLY prompt the model to output <external>	usage of audio beyond the provided splits is strictly	1073
1019	only when the internal prediction is wrong while	prohibited.	1074
	external perception is correct.		
1020		<b>D.2 Temporal Soundscapes QA</b>	1075
1021	You are an audio understanding model	This subset focuses on temporal reasoning over	1076
1022	with access to three inputs:	sound events, encompassing 26 event classes.	1077
1023	(1) The original audio;	Questions require identifying active sound classes,	1078
1024	(2) One answer candidate generated by	temporal ordering, timestamp estimation (onset,	1079
1025	another model (external);	offset, duration), and event comparison. The sub-	1080
1026	(3) Your own prediction (internal).	set comprises approximately 1.0K training and	1081
1027		0.6K development QA pairs. Audio clips are	1082
1028	Your task is to decide which of the	mono-channel, 10 seconds long, and sampled at	1083
1029	following strategies to apply:	32–48 kHz. Most clips correspond to a single	1084
1030	- If your internal prediction is correct	QA item, while a small portion supports multiple	1085
1031	and acceptable, output <internal> and	questions. All annotations are manually verified,	1086
1032	repeat your answer.	which include event types, timestamps, and an-	1087
1033	- If the external candidate is correct	swers. Audio is sourced from NIGENS, L3DAS23	1088
1034	while your internal prediction is	Challenge (Marinoni et al., 2024), and TAU Spa-	1089
1035	incorrect, output <external> and use	tial Sound Events 2019 datasets (Adavanne et al.,	1090
1036	the external answer.	2019).	1091
1037	- If all given answers are incorrect,		
1038	output <rewrite> and re-answer the	<b>D.3 Complex QA (MMAU)</b>	1092
1039	question correctly based only on the	This subset evaluates high-level reasoning over nat-	1093
1040	original audio.	ural sound scenes. Each instance consists of a 10-	1094
1041		second, 16 kHz audio clip paired with a question	1095
1042	Return the selected token	requiring reasoning over acoustic, temporal, and	1096
1043	(<internal>/<external>/<rewrite>)	contextual cues ( <i>i.e.</i> , overlapping events, implied	1097
1044	followed by your final answer.	sequences, or abstract relationships.) Tasks are in-	1098
		spired by the MMAU Sound benchmark (Sakshi	1099
1045		et al., 2024b) and extend its principles to more di-	1100
1046	<b>C Dataset Details for ASR</b>	verse soundscapes. The data include approximately	1101
1047	To ensure a fair and balanced evaluation across	6.4K training and 1.6K development QA pairs. Au-	1102
1048	diverse speech corpora, we uniformly sampled up	dio clips are sourced from AudioSet and Mira.	1103
1049	to 20,000 audio-question pairs from each dataset		
1050	for training as in Table 7. All datasets were aligned	<b>E Cases On ASR</b>	1104
1051	to a consistent prompt-question-answer format to	We provide several examples to show the effective-	1105
1052	support unified multi-dataset training.	ness of our Speech-Hands in ASR tasks.	1106
1053			
1054	<b>D Dataset Details for DCASE2025</b>		
1055	<b>AudioQA</b>		
1056	The DCASE2025 AudioQA benchmark comprises		
1057	three complementary multiple-choice question-		
1058	answering subsets, each designed to evaluate a dif-		
1059	ferent aspect of audio reasoning.		
	<b>D.1 Bioacoustics QA</b>		
	This subset targets perceptual and cognitive ground-		
	ing in marine bioacoustics. It includes questions		

Dataset	AMI	Tedlium	gigaspeech	spgispeech	Voxpopuli	Libri	Libri-clean	Libri-other
<b>Sampling #</b>	Subset in Speech-Hands Training							
<b>Train</b>	20,000	20,000	9,389	20,000	20,000	20,000	20,000	20,000

Table 7: Number of training samples used from each dataset. For GigaSpeech, only 9,389 valid samples met our filtering criteria.

Subset	#Train / #Dev	<in>/<ex>/<re>
Bio-acoustic QA	0.7K / 0.2K	338/234/168
Soundscape QA	1.0K / 0.6K	604/182/252
Complex QA	6.4K / 1.6K	4,267/1,785/391

Table 8: Statistics of the three DCASE2025 AudioQA subsets used in our experiments. <in>/<ex>/<re> shows the token distribution in training w/ majority sampling.

#### ASR Case 1

**Q.** What is the correct transcription of this audio?

**Internal:** water and fresh stalls to take on there

**External:** we’ve water and fresh stores to tick on there

**Rewrite:** we’ve water and fresh stores to take on there

**Pred:** <rewrite> we’ve water and fresh stores to take on there

**Gold:** we’ve water and fresh stores to take on there

#### ASR Case 2

**Q.** What is the correct transcription of this audio?

**Internal:** then has nothing happened with a voice like a whip mr wicker said start at once

**External:** then as nothing happened with a voice like a whip mister wigger said start at once

**Rewrite:** then as nothing happened with a voice like a whip mister wuthers said start at once

**Pred:** <external> then as nothing happened with a voice like a whip mister wigger said start at once

**Gold:** then as nothing happened with a voice like a whip mister wicker said start at once

## F Cases On AudioQA

For each token we pick up 5 examples to illustrate how action tokens influence the final correctness. For <rewrite> token, it works not only when the internal prediction is conflicted with the external prediction, but also show robust and accurate decision when hypotheses generate the same but wrong prediction. This is extremely important when facing challenging tasks, the omni-agent should be confident to refuse wrong hypotheses. Also for <internal> and <external> tokens that have more frequencies, they should be accurately distinguished when conflicts, and action tokens complete a good job on recognizing the true answer.

#### AudioQA Case 1

**Q.** Based on the audio, which natural phenomenon could be occurring?

*Options:*

- A. Earthquake
- B. Thunderstorm
- C. Forest fire
- D. Snowstorm

**Internal pred:** B. Thunderstorm

**External pred:** B. Thunderstorm

**Pred:** <rewrite> C. Forest fire

**Gold:** <rewrite> C. Forest fire

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1107

1108

### AudioQA Case 2

**Q.** What activity is likely occurring based on the described sounds?

*Options:*

- A. Driving a car
- B. Planting a garden
- C. Painting a picture
- D. Cooking a meal

**Internal pred:** A. Driving a car

**External pred:** B. Planting a garden

**Pred:** <rewrite> D. Cooking a meal

**Gold:** <rewrite> D. Cooking a meal

1126

### AudioQA Case 3

**Q.** Based on the audio, what is the most likely environment in which the sounds are occurring?

*Options:*

- A. An indoor setting
- B. A rural countryside
- C. A busy city street
- D. A quiet residential neighborhood

**Internal pred:** A. An indoor setting

**External pred:** B. A rural countryside

**Pred:** <rewrite> C. A busy city street

**Gold:** <rewrite> C. A busy city street

1127

### AudioQA Case 4

**Q.** What background sound is present in the audio clip?

*Options:*

- A. Waves crashing
- B. Birds chirping
- C. Bells ringing
- D. Running car engine

**Internal pred:** A. Waves crashing

**External pred:** D. Running car engine

**Pred:** <rewrite> C. Bells ringing

**Gold:** <rewrite> C. Bells ringing

1128

### AudioQA Case 5

**Q.** What might the purpose of tapping the metal object in the background be?

*Options:*

1129

- A. To emphasize the speaker's instructions
- B. To distract from the speaker's sad tone
- C. To demonstrate its material quality
- D. To create a rhythmic background

**Internal pred:** A. To emphasize the speaker's instructions

**External pred:** A. To emphasize the speaker's instructions

**Pred:** <rewrite> C. To demonstrate its material quality

**Gold:** <rewrite> C. To demonstrate its material quality

1130

### AudioQA Case 6

**Q.** What type of mood is conveyed through the musical elements in this audio?

*Options:*

- A. Calm and soothing
- B. Angry and aggressive
- C. Joyful and uplifting
- D. Sad and reflective

**Internal pred:** B. Angry and aggressive

**External pred:** D. Sad and reflective

**Pred:** <external> D. Sad and reflective

**Gold:** <external> D. Sad and reflective

1131

### AudioQA Case 7

**Q.** Based on the audio description, what is likely happening in the background?

*Options:*

- A. A calm evening
- B. A windy day
- C. An earthquake
- D. A quiet morning

**Internal pred:** B. A windy day

**External pred:** C. An earthquake

**Pred:** <external> C. An earthquake

**Gold:** <external> C. An earthquake

1132

### AudioQA Case 8

**Q.** What type of sound is present in the background of the audio clip?

*Options:*

- A. Car engine

1133

B. Ocean waves  
C. Upbeat synthesized music  
D. Bird chirping  
**Internal pred:** B. Ocean waves  
**External pred:** C. Upbeat synthesized music

**Pred:** <external> C. Upbeat synthesized music  
**Gold:** <external> C. Upbeat synthesized music

#### AudioQA Case 9

**Q.** Based on the audio description, what is the primary focus of the sounds?  
*Options:*  
A. A quiet library setting  
B. A peaceful nature scene  
C. A busy city street  
D. A defense attack scenario

**Internal pred:** A. A quiet library setting  
**External pred:** D. A defense attack scenario

**Pred:** <external> D. A defense attack scenario  
**Gold:** <external> D. A defense attack scenario

#### AudioQA Case 10

**Q.** Based on the audio description, what type of activity is most likely taking place?  
*Options:*  
A. Gardening  
B. Woodworking  
C. Lock-picking  
D. Cooking

**Internal pred:** B. Woodworking  
**External pred:** C. Lock-picking

**Pred:** <external> C. Lock-picking  
**Gold:** <external> C. Lock-picking

#### AudioQA Case 11

**Q.** What element in the audio contributes to the emotional depth of the song besides the vocals?  
*Options:*  
A. The language spoken  
B. The steady drum beats

C. The groovy bass line  
D. The keyboard harmony  
**Internal pred:** D. The keyboard harmony  
**External pred:** C. The groovy bass line

**Pred:** <internal> D. The keyboard harmony  
**Gold:** <internal> D. The keyboard harmony

#### AudioQA Case 12

**Q.** Based on the audio description, what type of environment is suggested by the background sounds?  
*Options:*  
A. A serene forest  
B. A quiet library  
C. A beach with waves  
D. A busy city street

**Internal pred:** D. A busy city street  
**External pred:** D. A busy city street

**Pred:** <internal> D. A busy city street  
**Gold:** <internal> D. A busy city street

#### AudioQA Case 13

**Q.** What might the purpose of the whistle and crinkling leaves in the background be?  
*Options:*  
A. To create a suspenseful atmosphere  
B. To signal the attention of someone nearby  
C. To mimic a bustling city environment  
D. To indicate the presence of wildlife

**Internal pred:** B. To signal the attention of someone nearby  
**External pred:** B

**Pred:** <internal> B. To signal the attention of someone nearby  
**Gold:** <internal> B. To signal the attention of someone nearby

#### AudioQA Case 14

**Q.** Why does the conversation feature electronic beats and rhythmic cymbal sounds?  
*Options:*  
A. To mimic the sounds of a busy environment

1134

1135

1136

1137

1138

1139

1140

1141

- B. To create a sense of urgency in the interaction
- C. To drown out background noise
- D. To enhance the emotional depth of the speech

**Internal pred:** B. To create a sense of urgency in the interaction

**External pred:** D. To enhance the emotional depth of the speech

**Pred:** <internal> B. To create a sense of urgency in the interaction

**Gold:** <internal> B. To create a sense of urgency in the interaction

1142

#### AudioQA Case 15

**Q.** What sound appears earliest in the audio?

*Options:*

- A. Ticking
- B. Accelerating, revving, vroom
- C. Idling
- D. Car

**Internal pred:** C. Idling

**External pred:** D. Car

**Pred:** <internal> C. Idling

**Gold:** <internal> C. Idling

1143