

# ChronosAudio: A Comprehensive Long-Audio Benchmark for Evaluating Audio-Large Language Models

Anonymous ACL submission

## Abstract

Although Audio Large Language Models (ALLMs) have witnessed substantial advancements, their long audio understanding capabilities remain unexplored. A plethora of benchmarks have been proposed for general audio tasks, they predominantly focus on short-form clips, leaving without a consensus on evaluating ALLMs over extended durations. This paper proposes *ChronosAudio*, the first multi-task benchmark tailored for long-audio understanding in ALLMs. It encompasses six major task categories and comprises 36,000 test instances totaling over 200 hours audio, stratified into short, middle, and long-form categories to comprehensively evaluate length generalization. Extensive experiments on 16 state-of-the-art models using *ChronosAudio* yield three critical findings: ① **Precipitous Long-Context Collapse**: ALLMs exhibit a severe inability to sustain performance, with the transition from short to long contexts triggering a staggering performance degradation of over 90% in specific tasks. ② **Structural Attention Dilution**: Performance degradation stems from a fundamental failure in maintaining temporal locality; attention mechanisms suffer from significant diffusion in later sequences. ③ **Restorative Ceiling of Mitigation**: Current strategies only offer 50% recovery. These findings reveal significant challenges in long-audio, underscoring the urgent need for approaches to achieve robust, document-level audio reasoning. The codes and dataset are available at <https://anonymous.4open.science/r/ChronosAudio-D49A>

## 1 Introduction

Recent advancements in Large Language Models (LLMs) have catalyzed the evolution of Multimodal LLMs (MLLMs), significantly enhancing their capabilities in audio understanding (Wang et al., 2025b), instruction following (Feng et al., 2025), and reasoning (Lu et al., 2025). Build-

Benchmark	Multi Dim.	Long Task	Total Dur.	Avg. Dur.
AudioBench (Wang et al., 2023)	✓	✗	400h	14s
AIR-Bench (Yang et al., 2024)	✓	✗	252h	35s
MuChoMusic (Weck et al., 2024)	✓	✗	5h	26s
MMSU (Wang et al., 2025a)	✓	✗	10h	7s
MMAU (Sakshi et al., 2025)	✓	✗	28h	10s
MMAR (Ma et al., 2025)	✓	✗	6h	19s
BLAB (Ahia et al., 2025)	✗	✓	833h	51m
<b>ChronosAudio (Ours)</b>	✓	✓	237h	322s

Table 1: Comparison between *ChronosAudio* with other audio benchmarks.

ing upon these foundations, state-of-the-art Audio Large Language Models (ALLMs) (An et al., 2024; Fan et al., 2025; Wu et al., 2025; Ding et al., 2025) have emerged to handle various audio-related tasks (Huang et al., 2023; Sun et al., 2025; Tian et al., 2025; Chen et al., 2025; Ghosh et al., 2025).

The burgeoning field of ALLMs has spurred the emergence of numerous datasets (Li et al., 2025a; Cheng et al., 2025; Lin et al., 2025; Li et al., 2025d) aimed at broadly evaluating their general capabilities. Representative benchmarks include Dynamic-SUPERB (Tseng et al., 2023), Audio-Bench (Wang et al., 2023), and AIR-Bench (Yang et al., 2024), alongside those prioritizing safety and trustworthiness like AudioTrust (Li et al., 2025c) and Audio Jailbreak (Song et al., 2025). However, Table 1 shows that their data is confined to short audio (Li et al., 2025d; Zhang et al., 2025), typically spanning merely seconds to a few minutes. Now systematic research into the long-duration audio understanding and analysis capabilities of ALLMs remains largely unexplored. Consequently, this raises a natural question: *How do Audio Large Language Models perform on long audio?*

To explore this question systematically, we introduce *ChronosAudio*, the first multi-task benchmark tailored for long-context evaluating (Ding et al., 2024; Li et al., 2024; Gao et al., 2025) on ALLMs as shown in Figure 1. *ChronosAudio*



Figure 1: A framework of ChronosAudio Benchmark for evaluating ALLMs across six major task categories.

is composed of **six major task categories, covering key long-audio application scenarios including Dictation, Localization, Transcription, Multi-Speaker, Comprehension, and Summary.** In addition, ChronosAudio contains **3,000** short-form audios (30s~5min), **2,000** middle-form audios (5min~10min) and **1,000** long-form audios (10min~20min) to provide a more comprehensive evaluation of the ALLMs' capabilities on long contexts. ChronosAudio includes **36,000** test instances, over **200 h** audio, and an average length of **322 s**.

Building on the *ChronosAudio* benchmark, our extensive experimentation rigorously demonstrates that state-of-the-art ALLMs **exhibit a critical vulnerability in maintaining high-fidelity understanding when processing extended temporal contexts.** For instance, we observe a precipitous performance collapse across evaluated models, where the transition from short to long contexts triggers a staggering degradation of over 90% in specific tasks. Even top-tier proprietary models are not immune to this decline; while closed-

source models achieve a Transcription score of 42.66 on short audio, this metric plummets to a mere 3.86 when processing long-form sequences. However, a distinct performance chasm is observed between model categories: while closed-source models maintain robust reasoning capabilities with a Summary score of 59.22 on long audio, open-source counterparts suffer a structural collapse to an average score of 14.14 in the same scenarios.

### Experimental Takeaways.

- **Precipitous Long-Context Collapse.** ALLMs exhibit a severe inability to sustain performance over extended durations. Our evaluation reveals that the transition from short to long contexts triggers a staggering performance degradation of over **90%** in specific tasks, with the *Lost-in-the-Middle* phenomenon manifesting significantly earlier than in other textual models.
- **Structural Attention Dilution.** Performance degradation stems from a fundamental failure in maintaining temporal locality. Visualization of attention weights confirms that while models

capture local dependencies in early stages, the attention mechanism suffers from significant diffusion in later sequences, resulting in a loss of high-resolution alignment.

- **Restorative Ceiling of Mitigation.** Current long-context strategies offer limited recovery. While mechanisms like Sparse Attention improve stability, they encounter a ceiling, recovering only **50%** of the model’s short-context proficiency, highlighting simple architectural tweaks are insufficient for document-level audio reasoning.

## 2 Preliminary & Related Work

### 2.1 Audio Large Language Models

Recent ALLMs have demonstrated significant potential in unifying tasks such as automatic speech recognition (Min and Wang, 2023; Bai et al., 2024a; Xie et al., 2025; Carone et al., 2025), speech translation (Huang et al., 2023; Liu et al., 2025), audio generation (Wang et al., 2025c), and dialogue (Zhang et al., 2023; Gong et al., 2024).

Prevailing ALLMs typically employ a two stage paradigm that aligns continuous audio signals with the discrete textual space of LLMs. In the first stage, the continuous audio waveform  $x_a \in \mathbb{R}^{T_a}$  is discretized into a sequence of acoustic tokens to bridge the modality gap. This is generally achieved via a neural audio tokenizer or codec  $\phi_a$ , which maps the input to a discrete codebook space  $\mathbb{Z}^{L_a}$ :

$$\mathbf{c}_a = \phi_a(x_a), \quad \phi_a : \mathbb{R}^{T_a} \rightarrow \mathbb{Z}^{L_a}, \quad (1)$$

where  $L_a$  denotes the length of result tokens.

To enable multimodal interaction, these acoustic tokens are projected into a shared semantic embedding space alongside textual inputs. Specifically, given a textual prompt  $x_t \in \mathbb{Z}^{L_t}$ , the model constructs a unified input sequence  $\mathbf{z}$  by concatenating the embeddings of both modalities:

$$\mathbf{z} = [\phi_e(\mathbf{c}_a) \parallel \phi_t(x_t)] \in \mathbb{R}^{(L_a+L_t) \times d}, \quad (2)$$

where  $\phi_e$  and  $\phi_t$  serve as the embedding functions and  $\parallel$  denotes the concatenation operation.

The unified embedding sequence is subsequently processed by a Transformer decoder  $f_\theta$ , which models the probability distribution over a joint vocabulary  $\mathcal{V}$  containing textual and acoustic units.  $\mathbf{h} = f_\theta(\mathbf{z})$  are projected via a matrix  $W \in \mathbb{R}^{|\mathcal{V}| \times d}$  to compute the next-token probabilities:

$$\hat{y} = \text{softmax}(W\mathbf{h}). \quad (3)$$

Eq 3 often referred to as joint autoregressive decoding, allows ALLMs to coherently generate interleaved audio and textual outputs, thereby facilitating complex tasks such as audio captioning (Guo et al., 2025) and dialogue (Gosai et al., 2025).

### 2.2 Long Context Understanding

In the realm of text-modal LLMs, long-context understanding, defined as the capacity to process and reason over extended sequences, has witnessed systematic advancements. Recent work has examined how LLMs process very long contexts (Liu et al., 2024; Wang et al., 2024; Kuratov et al., 2024; Wu et al., 2024; Yang et al., 2025). These studies include stress tests for retrieval in book-length documents, where performance degrades when key information is located in the middle part (Bai et al., 2024b,c; Que et al., 2024). These works have established a solid foundation for evaluating retrieval and reasoning in long-text scenarios.

Despite advancements in NLP, research on long-context understanding in the audio modality remains scarce, primarily because audio is continuous, making reasoning over long sequential structures, significantly more challenging than with discrete text tokens (Zhu et al., 2025). Furthermore, current sota ALLMs lack sufficient benchmarking for these long scenarios, as existing evaluations predominantly focus on short clips (Ahia et al., 2025; Ma et al., 2025; Wang et al., 2025a; Sakshi et al., 2025; Li et al., 2025b). This highlights a critical need for dedicated benchmarks and modeling approaches that target long-context understanding.

## 3 ChronosAudio

We introduce *ChronosAudio*, a multi-task benchmark specifically designed to assess long-audio understanding on **16** sota ALLMs. The benchmark organizes six tasks into three distinct categories: *Temporal Perception*, *Verbatim Sequential Generation*, and *High-level Reasoning*.

### 3.1 An Overview of ChronosAudio

To ensure a comprehensive evaluation of length generalization and address the scarcity of long-context benchmarks, *ChronosAudio* explicitly focuses on time-extended audio scenarios. The benchmark aggregates a diverse collection of audio data, totally amounting to over **200** hours.

Specifically, ChronosAudio comprises **36,000** test instances derived from **6,000** distinct audio

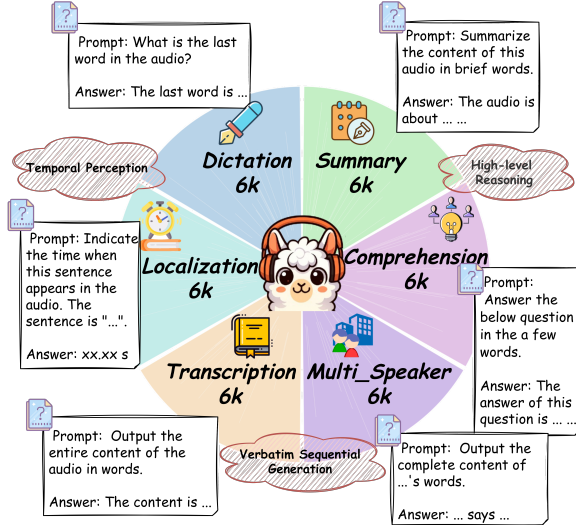


Figure 2: Overview of the ChronosAudio benchmark data composition. Each task comprises 6,000 test instances, totaling 36,000 samples.

clips. These instances span six major task categories, ensuring broad coverage of key long-context application scenarios. The audio data is rigorously stratified into three duration categories to probe model capabilities across different context windows. We divide data into three types:

- **Short-form:** 3,000 audio clips (30s~5min), serving as a baseline for context handling.
- **Middle-form:** 2,000 audio clips (5min~10min), bridging the gap between sentence-level and document-level processing.
- **Long-form:** 1,000 audio clips (10min~20min), specifically designed to challenge the limits of long-context attention and reasoning.

Figure 2 shows that all datasets in ChronosAudio are standardized into a unified format. This standardization is critical to facilitate consistent, automated evaluation across diverse audio-centric capabilities and to enable fair comparisons between different performance of ALLMs.

### 3.2 Temporal Perception

This category focuses on the model’s ability to locate specific information in the time domain and sustain attention over long durations.

**Dictation.** This task probes attention persistence. The metric, Dictation Accuracy ( $A_{\text{dict}}$ ), is calculated as the ratio of correctly transcribed units ( $N_{\text{correct}}$ ) to total reference units ( $N_{\text{total}}$ ):

$$DS = \frac{N_{\text{correct}}}{N_{\text{total}}}, \quad (4)$$

where  $N_{\text{correct}}$  denotes the count of units (words or

characters) correctly transcribed by the model, and  $N_{\text{total}}$  represents the total number of units in the reference text.

**Localization.** This task measures temporal grounding. We first compute the absolute temporal error  $\Delta t = |t_{\text{pred}} - t_{\text{ref}}|$ . With a strict tolerance threshold  $\tau = 0.1s$ , the localization score is defined as:

$$LS = \begin{cases} 1 - \frac{\Delta t}{\tau}, & \text{if } \Delta t \leq \tau \\ 0, & \text{if } \Delta t > \tau \end{cases} \quad (5)$$

where  $t_{\text{pred}}$  and  $t_{\text{ref}}$  represent the predicted timestamp and the ground-truth timestamp, respectively.  $\tau$  serves as the tolerance boundary, penalizing predictions that deviate beyond 0.1 seconds.

### 3.3 Verbatim Sequential Generation

This category evaluates the model’s capacity to generate long, faithful transcripts, preserving exact lexical content and speaker attribution.

**Transcription.** We evaluate transcription using a composite metric reflecting word accuracy and semantic similarity. The final score is:

$$TS = \beta_1(1 - \text{mWER}) + \beta_2 \text{BertSc}, \quad (6)$$

where mWER denotes the mixed Word Error Rate, and BERTScore measures semantic similarity.  $\beta_1$  and  $\beta_2$  are the balancing coefficients.

**Multi-Speaker.** This task requires consistent speaker IDs and precise turn boundaries. Evaluation combines content accuracy with:

$$MS = \gamma_1(1 - \text{mWER}) + \gamma_2 \left( \frac{2P_{\text{spk}}R_{\text{spk}}}{P_{\text{spk}} + R_{\text{spk}}} \right), \quad (7)$$

where  $P_{\text{spk}}$  and  $R_{\text{spk}}$  denote the precision and recall of speaker attribution, respectively. The weights are configured as  $\gamma_1$  and  $\gamma_2$  to balance the content accuracy and speaker identification.

### 3.4 High-level Reasoning

This category assesses the model’s ability to synthesize information and perform logical reasoning. We define  $\mathcal{I}(\cdot)$  as the information extraction function and  $\mathcal{K}(\cdot)$  as the key point mapping function.

**Comprehension.** We evaluate the QA capability based on exact match and a hallucination penalty. The final score integrates these two factors:

$$CS = \left[ w_1 \cdot \mathbb{I}(\text{ExactMatch}) - w_2 \cdot \frac{|\mathcal{I}(\mathcal{Y}) \setminus \mathcal{I}(\mathcal{G})|}{|\mathcal{I}(\mathcal{Y})|} \right], \quad (8)$$

Model	Dictation			Localization			Transcription			Multi-Speaker			Comprehension			Summary			
	S	M	L	S	M	L	S	M	L	S	M	L	S	M	L	S	M	L	
<b>Open Source Models</b>																			
Qwen2-Audio-7B	15.30	7.90	0.00	4.40	0.00	0.00	25.80	10.00	0.00	16.60	<b>9.50</b>	<b>5.00</b>	5.00	27.90	6.60	0.00	28.00	25.60	9.80
Qwen2.5-Omni-3B	22.50	14.90	0.00	6.00	0.00	0.00	50.10	16.80	0.00	15.70	7.80	3.00		30.50	18.40	<b>10.20</b>	31.10	25.80	4.60
Qwen2.5-Omni-7B	24.80	<b>15.40</b>	0.00	6.70	0.00	0.00	<b>56.10</b>	18.70	0.00	17.20	9.60	4.10		<b>34.00</b>	<b>20.20</b>	<b>11.30</b>	34.20	28.70	5.60
Qwen3-Omni-30B	<b>27.90</b>	<b>15.00</b>	0.00	<b>7.10</b>	0.00	0.00	53.90	17.20	0.00	19.00	8.10	3.90		31.80	<b>19.00</b>	9.70	77.50	56.00	18.90
MiniCPM-o-2_6	9.30	2.60	0.00	3.70	0.00	0.00	10.50	1.40	0.00	4.90	0.00	0.00		15.80	9.00	0.00	54.60	39.00	21.60
Kimi-Audio-7B	0.00	0.00	0.00	0.00	0.00	0.00	5.00	0.90	0.00	0.90	0.00	0.00		6.70	0.00	0.00	33.80	27.00	13.40
SALMONN	0.00	0.00	0.00	0.00	0.00	0.00	9.80	3.30	0.00	1.80	0.00	0.00		11.20	0.00	0.00	40.90	38.40	30.20
Step-Audio-2-mini	10.50	3.80	0.00	0.00	0.00	0.00	13.30	7.60	0.00	1.20	0.00	0.00		18.90	6.30	0.00	27.60	21.90	9.00
<b>Average</b>	<b>13.79</b>	<b>7.45</b>	<b>0.00</b>	<b>3.49</b>	<b>0.00</b>	<b>0.00</b>	<b>28.06</b>	<b>9.49</b>	<b>0.00</b>	<b>9.66</b>	<b>4.38</b>	<b>2.00</b>		<b>22.10</b>	<b>10.65</b>	<b>3.90</b>	<b>40.96</b>	<b>32.80</b>	<b>14.14</b>
<b>Closed Source Models</b>																			
GPT-4o-Audio	9.10	2.30	0.00	5.80	0.00	0.00	41.90	16.50	5.60	13.90	7.80	2.80		28.30	12.00	5.70	78.70	72.10	<b>73.30</b>
Gemini-3.0-Pro	<b>26.70</b>	5.00	<b>0.30</b>	<b>7.40</b>	0.00	0.00	45.80	<b>24.20</b>	<b>13.90</b>	11.30	7.80	0.50		<b>32.00</b>	13.90	4.30	<b>81.20</b>	<b>75.40</b>	70.30
Gemini-2.5-Pro	12.10	2.90	0.00	6.10	0.00	0.00	37.00	<b>22.90</b>	10.10	12.50	7.30	0.40		31.20	9.20	3.40	<b>79.50</b>	<b>72.30</b>	<b>70.40</b>
Gemini-2.5-Flash	8.80	1.50	0.20	2.80	0.00	0.00	49.30	17.30	9.40	9.50	5.30	1.00		30.50	11.90	5.50	73.70	64.50	26.30
Gemini-2.5-Flash-Lite	9.30	1.50	0.10	3.70	0.00	0.00	41.60	14.00	1.40	<b>20.60</b>	<b>10.10</b>	<b>4.30</b>		26.50	11.50	4.90	71.00	50.80	62.30
Gemini-2.0-Flash	11.50	2.50	0.00	1.90	0.00	0.00	38.80	14.10	0.00	8.00	2.40	0.00		23.00	14.20	0.00	76.20	63.00	68.10
Gemini-2.0-Flash-Lite	9.00	2.40	0.00	2.00	0.00	0.00	41.70	16.30	2.90	7.20	1.00	0.00		21.50	13.60	0.00	75.40	66.30	66.10
Qwen-Omni-Turbo	25.00	13.60	<b>1.20</b>	6.00	0.00	0.00	<b>58.70</b>	20.10	<b>11.90</b>	<b>19.90</b>	9.30	0.00		27.80	18.30	7.30	76.40	62.60	60.50
<b>Average</b>	<b>9.54</b>	<b>2.04</b>	<b>0.06</b>	<b>3.24</b>	<b>0.00</b>	<b>0.00</b>	<b>42.66</b>	<b>15.64</b>	<b>3.86</b>	<b>11.84</b>	<b>5.32</b>	<b>1.62</b>		<b>25.96</b>	<b>12.64</b>	<b>3.89</b>	<b>75.00</b>	<b>63.34</b>	<b>59.22</b>

Table 2: Performance comparison on six tasks with different audio lengths (Short/Middle/Long). The column order aligns with the visual analysis charts, averaged over 5 experimental rounds. The best and second-best performing models are highlighted in **bold** with 🏆 and 🥈 medals respectively.

where  $\mathcal{Y}$  and  $\mathcal{G}$  denote the model prediction and ground truth, respectively.  $\mathbb{I}(\cdot)$  is the indicator function. The two terms correspond to Exact Match and Hallucination Rate, weighted by  $w_1$  and  $w_2$ .

**Summary.** Summarization quality is determined by Coverage and Factualty. Metric is defined as:

$$SS = \left[ \delta_1 \frac{|\mathcal{K}(\mathcal{Y}) \cap \mathcal{K}(\mathcal{G})|}{|\mathcal{K}(\mathcal{G})|} + \delta_2 \frac{|\mathcal{K}(\mathcal{Y}) \cap \mathcal{K}(\mathcal{G})|}{|\mathcal{K}(\mathcal{Y})|} \right]. \quad (9)$$

The terms represent Coverage (Recall) and Factualty (Precision), balanced by  $\delta_1$  and  $\delta_2$ .

## 4 Experiments and Evaluations

In this section, we conduct experiments to address the following research questions:

- *RQ1: How does the ALLMs performance evolve as the audio duration scales from short to long contexts?* Specifically, to what extent does performance degrade, and does a distinct decline emerge as the context length increases?
- *RQ2: What does attention analysis reveal about the underlying causes of performance degradation in extended audio sequences?* Through the visualization of attention weights, does the performance drop stem from a structural collapse or are there other internal factors at play?
- *RQ3: Can mitigation strategies effectively alleviate the performance collapse in entire long-*

*duration tasks?* In particular, are these mainstream mechanisms really capable of recovering all the lost capabilities within the subset?

### 4.1 Experimental Setup

**Baselines.** The baseline model pool for our evaluation is defined as Qwen2-Audio-7B (An et al., 2024), Qwen-Omni series (Fan et al., 2025; Xu et al., 2025), MiniCPM-o-2\_6 (Yao et al., 2024), Kimi-Audio-7B (Ding et al., 2025), SALMONN (Tang et al., 2023), and Step-Audio-2-mini (Wu et al., 2025). Furthermore, we include leading closed-source models to establish upper-bound performance benchmarks, including GPT-4o-Audio (Achiam et al., 2023), Gemini-Flash and Gemini-Pro series (Comanici et al., 2025), and Qwen-Omni-Turbo. We access ALLMs via their respective official APIs.

**Parameter & Model Configurations.** The backbone LLM for our proposed framework is consistently initialized with Qwen2-Audio-7B. For the data processing stage, we categorize samples into three distinct audio lengths, *Short*, *Middle*, and *Long*, to make sure comprehensive evaluation coverage across temporal scales.

### 4.2 RQ1: Performance on Long Audio Tasks

To systematically evaluate the performance of state-of-the-art ALLMs in terms of long-context robust-

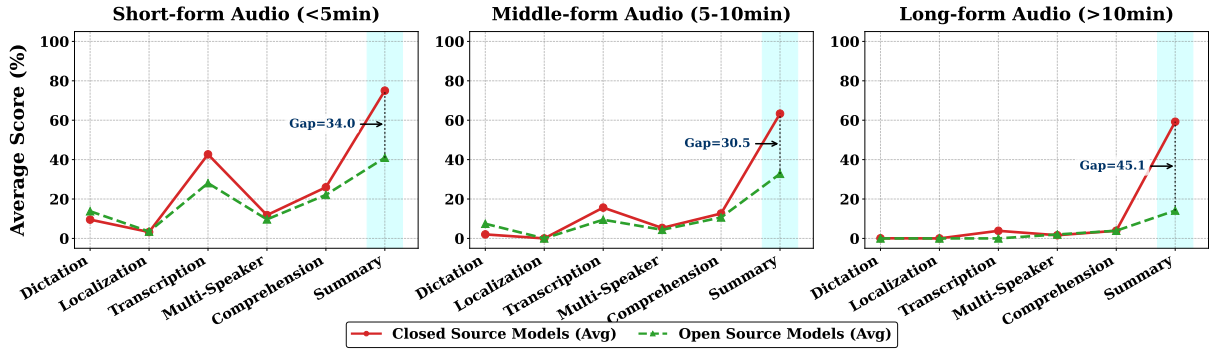


Figure 3: Average Performance between Closed-Source and Open-Source Models across Audio Durations.

ness, we conducted a comprehensive assessment on 16 models using the *ChronosAudio* benchmark. Table 2 presents the comparative results across 6 tasks, stratified by 3 audio duration. We can draw the following observations:

*Obs. 1: ALLMs exhibit a precipitous performance decline in long-context scenarios.*

Table 2 reveals that while most models achieve competitive performance on short audio, their capabilities degrade significantly on long audio, with a remarkable decline about 18 points, from 3.37 to 21.30  $\uparrow$ . (I) *Open-source degradation*: The average performance drop from Short to Long audio across all open-source baselines is substantial. Notably, in the Summary task, the average performance for open-source models plummets from 40.96 (Short) to 14.14 (Long), indicating a 26.82 collapse for extended durations. (II) *Closed-source limitations*: Similarly, even advanced closed-source models like the Gemini-2.0-Flash Series show a noticeable performance gap, whose synthesis score has dropped nearly 20 points. While closed-source models maintain a Transcription score of 42.66 on Short audio, this metric drops to 3.86 on Long audio. This highlights that modeling long-range acoustic dependencies remains a universal challenge, even for top-tier proprietary models.

*Obs. 2: Closed-source models exhibit superior long-context robustness and reasoning capabilities compared to open-source counterparts.*

Comparative analysis of Table 2 and Figure 3 reveals a distinct performance chasm between closed-source and open-source models, particularly in complex long-form tasks. While open-source models like Qwen2.5-Omni series achieve competitive parity in short-form perception, they lack stability in long-duration scenarios. For instance, the synthesis score for them has less 10 points than closed-source models. Particularly, the average score for

open-source models in the Summary (Long) task is merely 14.14, whereas closed-source models maintain a robust average of 59.22, outperforming them by a margin of over 45 points. The divergence is most acute in tasks requiring high-level information synthesis. In Comprehension (Long), open-source models suffer a near-total collapse (average score 3.90), while closed-source models retain significant utility with an average score 25.96.

*Obs. 3: The Lost-in-the-Middle phenomenon manifests as an early performance cliff, with significant degradation occurring immediately upon transitioning from Short to Middle durations.*

While models demonstrate robust capabilities on Short audio, Table 2 reveals a critical vulnerability: except Summary task and Localization task, the transition to middle-form triggers a disproportionate loss (over 12 points) in fidelity, rather than a gradual decline. This effect is most evident in the **Transcription** task. The average score for open-source models crashes from 28.06 (Short) to 9.49 (Middle), a staggering reduction of nearly 66%. Similarly, closed-source models are not immune, with their average performance more than halving from 42.66 to 15.64. This indicates that the Lost-in-the-Middle effect in ALLMs occurs much earlier than in text LLMs, suggesting that the effective context window for high-fidelity understanding is currently saturated at the short level.

### 4.3 RQ2: Analysis of Attention Mechanism

To investigate the stability of the attention mechanism in Qwen2-Audio-7B when processing long-duration sequences, we visualize the attention weights across three distinct temporal stages: **First**, **Middle**, and **Last** of the sequence.

*Obs. 4: The model fails to maintain precise temporal locality, exhibiting significant attention dilution as the sequence length increases.*

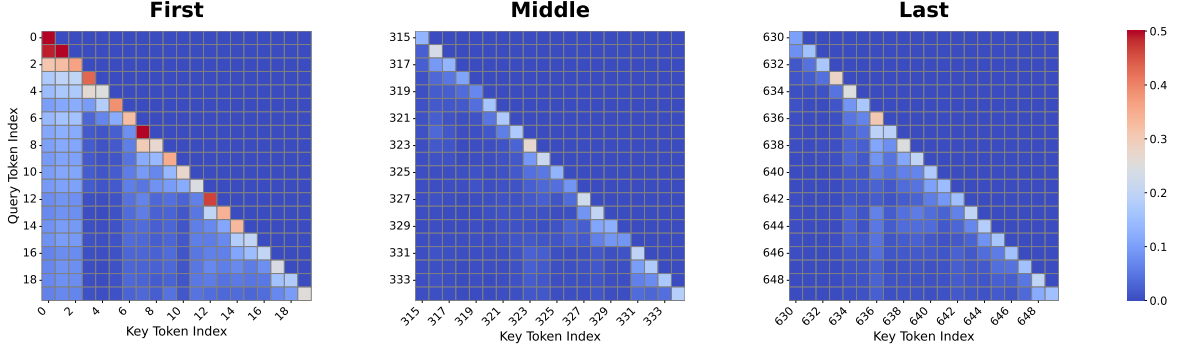


Figure 4: **Visualization of Attention Consistency.** The heatmaps display the self-attention weights for the First, Middle, and Last segments. Darker colors represent lower weights, while brighter colors indicate higher attention scores.

As illustrated in Figure 4, the attention heads exhibit sharp, distinct diagonal alignment in the **First** phase, indicating that the model effectively captures local dependencies with high precision at the beginning. However, this consistency is not sustained. In the **Middle** phase, the diagonal pattern begins to fade, showing early signs of degradation. By the **Last** phase, the diagonal alignment becomes significantly diffused and less defined compared to the First phase. This comparison reveals that the model struggles to maintain the same level of attention quality throughout the sequence, suffering from **structural attention dilution** where the high-resolution acoustic alignment is lost.

#### 4.4 RQ3: Efficacy of Mitigation Strategies

Based on analysis of RQ2, we evaluate the efficacy of two mainstream long-context mitigation strategies (Zaheer et al., 2020): **Sparse Attention** (Child et al., 2019) and **Sliding Window Attention** (Beltagy et al., 2020; Xiao et al., 2023), to address the performance collapse identified in RQ1.

**1) Sparse Attention (SA).** It retains a set of global tokens alongside the local window to preserve long-range dependencies. The attention score  $A_{ij}$  between query  $q_i$  and key  $k_j$  is defined as:

$$A_{ij}^{\text{Sparse}} = \frac{q_i k_j^\top}{\sqrt{d_k}}, \quad \text{where } j \in \mathcal{L}_i \cup \mathcal{G}, \quad (10)$$

where  $\mathcal{L}_i = \{j \mid i - w \leq j \leq i\}$  represents the local window,  $\mathcal{G}$  denotes the set of pre-selected global token indices and  $d_k$  is the head dimension. **2) Sliding Window Attention (SWA).** It restricts the receptive field of each token to a fixed-size local neighborhood, thereby reducing computational complexity. With  $w$  denotes the window

size, the  $A_{ij}$  formulation is:

$$A_{ij}^{\text{Window}} = \frac{q_i k_j^\top}{\sqrt{d_k}}, \quad \text{where } i - w \leq j \leq i. \quad (11)$$

*Obs. 5: Sparse Attention significantly revitalizes long-context performance, demonstrating superior efficacy over other mechanisms.*

As substantiated in Table 3, applying appropriate global attention constraints effectively mitigates the attention dilution inherent in extended sequences.

- SA consistently yields the most substantial recovery across all benchmarks. The most compelling evidence is observed in the **Transcription** task for Minicpm-o-2\_6. Here, Sparse Attention delivers a remarkable performance boost of  $\Delta = 14.4\%$ , elevating the score from a baseline of **17.20** to **31.60**. This suggests that by explicitly retaining global tokens, the model successfully isolates high-importance acoustic features essential for reasoning, thereby filtering out interference from irrelevant context.
- Conversely, while SWA reduces computational complexity, it imposes a rigid locality constraint that hinders global dependency modeling. This tunnel vision effect is evident in the **Transcription** task for Qwen2.5-Omni-3B, where SWA (**16.50**) lags distinctively behind SA (**25.20**).

*Obs. 6: Mitigation strategies exhibit a restorative ceiling, failing to fully bridge the performance gap between long-context tasks and short baselines.*

To assess the recovery from the Lost-in-the-Middle phenomenon, we benchmark the improved middle-form results against the models' intrinsic short-form capabilities, referring to Table 2.

- In specific retrieval-oriented tasks, Sparse Attention effectively eliminates the performance



## 6 Limitations

While *ChronosAudio* establishes a rigorous benchmark, our study acknowledges several limitations. For our exploration of mitigation strategies, although we identified a restorative ceiling where these methods recover only limited proficiency, we have yet to investigate training-based solutions, such as continued pre-training or frequency-based positional interpolation. Regarding data diversity, our benchmark is predominantly English-centric and from relatively clean acoustic environments. The cross-lingual generalization of ALLMs and their robustness against real-world noise remain critical avenues for future study. Furthermore, while we benchmarked a wide array of models, our microscopic analysis of Attention Dilution was restricted to open-source models due to the black-box nature of closed-source models.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Orevaoghene Ahia, Martijn Bartelds, Kabir Ahuja, Hila Gonen, Valentin Hofmann, Siddhant Arora, Shuyue Stella Li, Vishal Puttagunta, Mofetoluwa Adeyemi, Charishma Buchireddy, Ben Walls, and Noah Bennett. 2025. Blab: Brutally long audio bench. *arXiv preprint arXiv:2505.03054*.

Zheng An, Wenqiang Chen, Xiao Yang, Ziyao Song, Qizhen Liu, Xiaowei Feng, Kai Fu, Han Fang, Jian Wang, Guohai Xu, and 1 others. 2024. Qwen2-audio technical report. *arXiv preprint arXiv:2407.10759*.

Ye Bai, Jingping Chen, Jitong Chen, Wei Chen, Zhuo Chen, Chuang Ding, Linhao Dong, Qianqian Dong, Yujiao Du, Kepan Gao, and 1 others. 2024a. Seed-asr: Understanding diverse speech and contexts with llm-based speech recognition. *arXiv preprint arXiv:2407.04675*.

Yuntao Bai, Xin Lv, Jun Zhang, Haoran Lyu, Jie Tang, Zhiyuan Huang, Zhengxin Du, Xiaoyang Liu, Aohan Zeng, Lei Hou, and 1 others. 2024b. Longbench: A bilingual, multitask benchmark for long context understanding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3119–3137.

Yuntao Bai, Shilin Tu, Jun Zhang, Hao Peng, Xiaozhi Wang, Xin Lv, Shiyu Cao, Jinghan Xu, Lei Hou, Yuxiao Dong, and 1 others. 2024c. Longbench v2: Towards deeper understanding and reasoning on realistic long-context multitasks. *arXiv preprint arXiv:2412.15204*.

Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.

Brandon James Carone, Iran R Roman, and Pablo Ripollés. 2025. The muse benchmark: Probing music perception and auditory relational reasoning in audio llms. *arXiv preprint arXiv:2510.19055*.

Jingyi Chen, Zhimeng Guo, Jiyun Chun, Pichao Wang, Andrew Perrault, and Micha Elsner. 2025. Do audio llms really listen, or just transcribe? measuring lexical vs. acoustic emotion cues reliance. *arXiv preprint arXiv:2510.10444*.

Xize Cheng, Dongjie Fu, Chenyuhao Wen, Shannon Yu, Zehan Wang, Shengpeng Ji, Siddhant Arora, Tao Jin, Shinji Watanabe, and Zhou Zhao. 2025. Aha-bench: Benchmarking audio hallucinations in large audio-language models. In *Advances in Neural Information Processing Systems (NeurIPS)*.

Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. 2019. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*.

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.

Ding Ding, Zeqian Ju, Yichong Leng, Songxiang Liu, Tong Liu, Zeyu Shang, Kai Shen, Wei Song, Xu Tan, Heyi Tang, and 1 others. 2025. Kimi-audio technical report. *arXiv preprint arXiv:2504.18425*.

Yiran Ding, Li Lyna Zhang, Chengruidong Zhang, Yuanyuan Xu, Ning Shang, Jiahang Xu, Fan Yang, and Mao Yang. 2024. Longrope: Extending llm context window beyond 2 million tokens. *arXiv preprint arXiv:2402.13753*.

Yang Fan, Xuemei Li, Jiannan Xiang, Xiang Li, Shuzheng Si, Xiaoyang Wang, Zhipeng Chen, Wenjie Li, Jiangwei Li, Jianghao Li, and 1 others. 2025. Qwen2.5-omni technical report. *arXiv preprint arXiv:2503.20215*.

Kunyu Feng, Yue Ma, Xinhua Zhang, Boshi Liu, Yikuang Yuluo, Yinhan Zhang, Runtao Liu, Hongyu Liu, Zhiyuan Qin, Shanhui Mo, and 1 others. 2025. Follow-your-instruction: A comprehensive mllm agent for world data synthesis. *arXiv preprint arXiv:2508.05580*.

Chaochen Gao, Xing Wu, Zijia Lin, Debing Zhang, and Songlin Hu. 2025. Nextlong: Toward effective long-context training without long documents. *arXiv preprint arXiv:2501.12766*.

639	Sreyan Ghosh, Arushi Goel, Lasha Koroshinadze, Sanggil Lee, Zhifeng Kong, Joao Felipe Santos, Ramani Duraiswami, Dinesh Manocha, Wei Ping, Mohammad Shoeybi, and 1 others. 2025. Music flamingo: Scaling music understanding in audio language models. <i>arXiv preprint arXiv:2511.10289</i> .	693
640		694
641		695
642		696
643		697
644		698
645	Yuan Gong, Cheng-I Hong, Shang-Hua Chang, and James Glass. 2024. Listen, think, and understand. In <i>The Twelfth International Conference on Learning Representations</i> .	699
646		700
647		701
648		702
649	Advait Gosai, Tyler Vuong, Utkarsh Tyagi, Steven Li, Wenjia You, Miheer Bavare, Arda Uçar, Zhongwang Fang, Brian Jang, Bing Liu, and 1 others. 2025. Audio multichallenge: A multi-turn evaluation of spoken dialogue systems on natural human interaction. <i>arXiv preprint arXiv:2512.14865</i> .	703
650		704
651		705
652		706
653		707
654		708
655	Tianyu Guo, Hongyu Chen, Hao Liang, Meiyi Qiang, Bohan Zeng, Linzhuang Sun, Bin Cui, and Wentao Zhang. 2025. Brace: A benchmark for robust audio caption quality evaluation. <i>arXiv preprint arXiv:2512.10403</i> .	709
656		710
657		711
658		712
659		713
660	Zhichao Huang, Rong Ye, Tom Ko, Qianqian Dong, Shanbo Cheng, Mingxuan Wang, and Hang Li. 2023. Speech translation with large language models: An industrial practice. <i>arXiv preprint arXiv:2312.13585</i> .	714
661		715
662		716
663		717
664	Yuri Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rodkin, Dmitry Sorokin, Artyom Sorokin, and Mikhail Burtsev. 2024. Babilong: Testing the limits of llms with long context reasoning-in-a-haystack. In <i>Advances in Neural Information Processing Systems</i> .	718
665		719
666		720
667		721
668		722
669	Bohan Li, Wenbin Huang, Yuhang Qiu, Yiwei Guo, Hankun Wang, Zhihan Li, Jing Peng, Ziyang Ma, Xie Chen, and Kai Yu. 2025a. Isa-bench: Benchmarking instruction sensitivity for large audio language models. <i>arXiv preprint arXiv:2510.23558</i> .	723
670		724
671		725
672		726
673		727
674	Bohan Li, Wenbin Huang, Yuhang Qiu, Yiwei Guo, Hankun Wang, Zhihan Li, Jing Peng, Ziyang Ma, Xie Chen, and Kai Yu. 2025b. Isa-bench: Benchmarking instruction sensitivity for large audio language models. <i>arXiv preprint arXiv:2510.23558</i> .	728
675		729
676		730
677		731
678		732
679	Kai Li, Can Shen, Yile Liu, Jirui Han, Kelong Zheng, Xuechao Zou, Zhe Wang, Shun Zhang, Xingjian Du, Hanjun Luo, and 1 others. 2025c. Audiotrust: Benchmarking the multifaceted trustworthiness of audio large language models. <i>arXiv preprint arXiv:2505.16211</i> .	733
680		734
681		735
682		736
683		737
684		738
685	Tianyu Li, Guanzheng Zhang, Quang D Do, Xin Yue, and Weizhe Chen. 2024. Long-context llms struggle with long in-context learning. <i>arXiv preprint arXiv:2404.02060</i> .	739
686		740
687		741
688		742
689	Wenyu Li, Xiaoqi Jiao, Yi Chang, Guangyan Zhang, and Yiwen Guo. 2025d. Auditorole: An audio dataset for character role-playing in large language models. <i>arXiv preprint arXiv:2509.23435</i> .	743
690		744
691		745
692		746
		747
	Liang Lin, Miao Yu, Kaiwen Luo, Yibo Zhang, Lilan Peng, Dexian Wang, Xuehai Tang, Yuanhe Zhang, Xikang Yang, Zhenhong Zhou, Kun Wang, and Yang Liu. 2025. Hidden in the noise: Unveiling backdoors in audio llms alignment through latent acoustic pattern triggers. <i>arXiv preprint arXiv:2508.02175</i> .	
	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the middle: How language models use long contexts. <i>Transactions of the Association for Computational Linguistics</i> , 12:157–173.	
	Wei Liu, Jiahong Li, Yiwen Shao, and Dong Yu. 2025. Tta: Transcribe, translate and alignment for cross-lingual speech representation. <i>arXiv preprint arXiv:2511.14410</i> .	
	Haolang Lu, Yilian Liu, Jingxin Xu, Guoshun Nan, Yuanlong Yu, Zhican Chen, and Kun Wang. 2025. Auditing meta-cognitive hallucinations in reasoning large language models. <i>arXiv preprint arXiv:2505.13143</i> .	
	Ziyang Ma and 1 others. 2025. Mmar: A challenging benchmark for deep reasoning in speech, audio, music, and their mix. <i>arXiv preprint arXiv:2505.13032</i> .	
	Zeping Min and Jinbo Wang. 2023. Exploring the integration of large language models into automatic speech recognition systems: An empirical study. In <i>International Conference on Neural Information Processing</i> , pages 69–84.	
	Haoyang Que, Feilong Duan, Liang He, Yutao Mou, Wei Zhou, Jun Liu, Weiqiang Rong, Zhiming M Wang, Jian Yang, Guannan Zhang, and 1 others. 2024. Hellobench: Evaluating long text generation capabilities of large language models. <i>arXiv preprint arXiv:2409.16191</i> .	
	S Sakshi, Utkarsh Tyagi, Sonal Kumar, Ashish Seth, Ramaneswaran Selvakumar, Oriol Nieto, Ramani Duraiswami, Sreyan Ghosh, and Dinesh Manocha. 2025. Mmau: A massive multi-task audio understanding and reasoning benchmark. In <i>International Conference on Learning Representations (ICLR)</i> .	
	Zirui Song, Qian Jiang, Mingxuan Cui, Mingzhe Li, Lang Gao, Zeyu Zhang, Zixiang Xu, Yanbo Wang, Chenxi Wang, Guangxian Ouyang, and 1 others. 2025. Audio jailbreak: An open comprehensive benchmark for jailbreaking large audio-language models. <i>arXiv preprint arXiv:2505.15406</i> .	
	Zhe Sun, Yujun Cai, Jiayu Yao, and Yiwei Wang. 2025. Spatial blind spot: Auditory motion perception deficits in audio llms. <i>arXiv preprint arXiv:2511.13273</i> .	
	Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, and Chao Zhang. 2023. Salmonn: Towards generic hearing abilities for large language models. <i>arXiv preprint arXiv:2310.13289</i> .	



## Appendix

### 1 Experiment Setup

Our experiments are conducted using the state-of-the-art audio-large language models (ALLMs). The implementation relies on the Hugging Face transformers ecosystem. To handle the computational complexity of long audio sequences efficiently, we incorporate flash-attn (Flash Attention 2) for accelerated attention mechanisms. Additionally, librosa is employed for audio waveform pre-processing. Detailed hardware specifications and specific library versions are listed in the below Table.

Item	Value
Hardware	4 NVIDIA A800 GPUs, 80GB memory per GPU
CPU	Intel(R) Xeon(R) Gold 6430 (2.1GHz, 32 cores)
Memory	1TB DDR4 RAM
OS	Ubuntu 22.04 LTS (Linux ubuntu 5.15.0-25-generic)
Framework	torch==2.6.0, transformers==4.51.3, trl==0.9.6
CUDA / cuDNN	CUDA 12.4, cuDNN 9.1.0
Python Version	Python 3.10

Table 4: Computing Infrastructure Settings

### 2 Task Description

**T1: Dictation** This component probes whether a model can sustain attention over long recordings and correctly output the final lexical token.

**T2: Localization** This component evaluates temporal grounding: the model must find in the recording where it occurs, handling trailing silences, overlaps, and possible multiple occurrences.

**T3: Transcription** This component assesses verbatim long-form transcription quality, including punctuation and robustness to silences, and fillers.

**T4: Multispeaker** This component evaluates content transcription together with speaker attribution, requiring consistent IDs and precise turn to the target speaker. The model is tasked with exclusively transcribing utterances belonging to a specific target speaker.

Task	Prompt Template
<b>T1: Dictation</b>	"Listen to the audio corresponding to the given text, and what is the last word spoken in the audio?"
<b>T2: Localization</b>	"Listen to the audio corresponding to the given text, and indicate the time when this sentence appears in the audio. The sentence is ..."
<b>T3: Transcription</b>	"Listen to the audio corresponding to the given text, and output the entire content of the audio in text form."
<b>T4: Multi-Speaker</b>	"Listen to the audio corresponding to the given text, and output the complete transcript of ...'s words."
<b>T5: Comprehension</b>	"Listen to the audio corresponding to the given text, and answer the question in the shortest possible sentence. The question is '...'"
<b>T6: Summary</b>	"Listen to the audio corresponding to the given text, and summarize the content of this audio in brief words."

Table 5: Overview of the prompt templates.

**T5: Comprehension** This component assesses whether a model can integrate dispersed acoustic evidence to answer content questions faithfully and concisely. The model must answer specific content-related questions by integrating dispersed acoustic evidence found throughout the recording.

**T6: Summary** This component tests faithful compression of long recordings into concise summaries with explicit evidence grounding. The model must distill extensive audio recordings into concise summaries that capture core themes.

### 3 Justification for LLM-as-a-Judge

Evaluating the output of ALLMs, particularly in tasks requiring high-level reasoning, presents significant challenges for traditional lexical metrics. Metrics that rely heavily on n-gram overlaps often fail to capture semantic equivalence, logical consistency, and factual accuracy. For instance, a generated summary may be semantically correct but lexically distinct from the ground truth, leading to an unfairly low score under such metrics. More-

## Uniform Prompt Template for Ensemble LLM-as-a-Judge

### [System Instructions]

You are an impartial and expert judge. Your task is to evaluate the quality of an audio model’s generated response. You will be provided with the **Audio Transcript (Ground Truth)**, the **User Question**, and the **Model’s Response**. Please evaluate the response independently based on the following criteria:

- **Accuracy & Faithfulness:** Does the response accurately reflect the information in the audio transcript without hallucination?
- **Completeness:** Does the response address all aspects of the user’s question?

### [Input Data]

#### 1. Audio Transcript (Ground Truth):

*{transcript}*

#### 2. User Question:

*{question}*

#### 3. Model Generated Response:

*{response}*

### [Evaluation Output]

Please output your judgment in the following format:

1. **Analysis:** Provide a concise analysis of the pros and cons of the response.
2. **Score:** Give a final integer score from 1 to 100 (where 100 is perfect), according to the formula of each task.

Figure 6: **The prompt template used for the Ensemble LLM-as-a-Judge framework.** To ensure a robust and unbiased evaluation, we employ three diverse state-of-the-art LLMs, **DeepSeek-V3.1-Terminus**, **Qwen3-VL-235B-A22B-Thinking**, and **Kimi-K2-Thinking** as independent judges. Each model receives this identical prompt to evaluate the generated response against the ground truth transcript. The final quality score is calculated as the arithmetic mean of the scores provided by the three judges:  $S_{final} = \frac{1}{3}(S_{DeepSeek} + S_{Qwen} + S_{Kimi})$ .

904 over, although human annotation is commonly re- 925  
905 garded as a more reliable alternative, it also suffers 926  
906 from inherent limitations, including high cost, lim- 927  
907 ited scalability, and variability across annotators. 928  
908 Subjective judgments and inconsistent annotation 929  
909 criteria can further introduce bias and noise, mak- 930  
910 ing it difficult to ensure reproducibility and fair 931  
911 comparison across different models and datasets. 932

912 To address the limitations of rigid lexical met- 933  
913 rics, we adopt a robust LLM-as-a-Judge approach. 934  
914 Instead of relying on a single proprietary model, we 935  
915 utilize a diverse panel of three state-of-the-art Large 936  
916 Language Models: **DeepSeek-V3.1-Terminus**, 937  
917 **Qwen3-VL-235B-A22B-Thinking**, and **Kimi-K2-** 938  
918 **Thinking**. While recent studies verify that ad- 939  
919 vanced LLMs correlate highly with human judg- 940  
920 ments, individual models may still exhibit specific 941  
921 stylistic biases or self-preference. By employing 942  
922 a multi-judge framework, we mitigate these bi- 943  
923 ases through a consensus-based scoring mecha- 944  
924 nism. Specifically, we provide each judge with

the input audio transcript (ground truth), the user 925  
question, and the model’s generated response. The 926  
judges independently assess the *faithfulness*, and 927  
*completeness* of the answer. The final quality score 928  
is derived from the arithmetic mean of the three 929  
evaluations, ensuring a more objective, scalable, 930  
and semantically reliable metric that transcends 931  
surface-level text matching. 932

## 4 Case Study 933

To comprehensively evaluate the capabilities of 934  
ALLMs in handling long-audio inputs, we propose 935  
*ChronosAudio*, a specialized benchmark designed 936  
to probe the temporal reasoning and semantic un- 937  
derstanding of current models. This benchmark 938  
is structured around six distinct tasks. All tasks’ 939  
cases are in following figures: 940

## Dictation Task (Short)

### Example: Qwen2-Audio-7B

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and what is the last word spoken in the audio?

**Label:** Techniques

**Answer (Raw Input Format):**

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 500+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
What is the last word spoken in the audio?<|im_end|>
<|im_start|>assistant
", "predict": "The last word spoken in the audio is 'techniques.'" }
```

## Dictation Task (Middle)

### Example: Qwen2-Audio-7B

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and what is the last word spoken in the audio?

**Label:** Seasons

**Answer①(Raw Input Format):**

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 1,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
What is the last word spoken in the audio?<|im_end|>
<|im_start|>assistant
", "predict": "Sorry I can't answer this question." }
```

**Answer②(Raw Input Format):**

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 1,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
What is the last word spoken in the audio?<|im_end|>
<|im_start|>assistant
", "predict": "beat" }
```

## Dictation Task (Long)

### Example: Qwen2-Audio-7B

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and what is the last word spoken in the audio?

**Label:** King

### Answer①(Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 2,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
What is the last word spoken in the audio?<|im_end|>
<|im_start|>assistant
", "predict": "Sorry I can't answer this question. Could you give me more information about
this audio?" }
```

### Answer②(Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 2,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
What is the last word spoken in the audio?<|im_end|>
<|im_start|>assistant
", "predict": "nervous." }
```

### Answer③(Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 2,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
What is the last word spoken in the audio?<|im_end|>
<|im_start|>assistant
", "predict": "speech speech speech speech speech speech speech speech speech speech speech speech
speech speech speech speech speech speech speech speech speech speech speech speech speech speech
speech speech speech speech speech speech speech speech speech speech speech speech speech speech
speech speech speech speech speech speech speech speech speech speech speech speech speech speech
speech speech speech speech speech speech speech speech speech speech speech speech speech speech
speech speech speech speech speech speech speech speech speech speech speech speech speech speech
speech speech speech speech speech speech speech speech speech speech speech speech speech speech
speech speech speech speech speech speech speech speech speech speech speech speech speech speech
speech speech speech speech speech speech speech speech speech speech speech speech speech speech
speech speech speech speech speech speech speech speech speech speech speech speech speech speech
speech speech speech speech speech speech speech speech speech speech speech speech speech speech
speech speech speech speech speech speech speech speech speech speech speech speech speech speech
speech speech speech speech speech speech speech speech speech speech speech speech speech speech
speech speech speech speech speech speech speech speech speech speech speech speech speech speech
speech ... .." }
```

## Localization Task (Short)

### Example: Qwen2-Audio-7B

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and indicate the time when this sentence appears in the audio. The sentence is "In this area the general water flow is distorted by more local conditions reaching across to France."

**Label:** 35.78 seconds

### Answer (Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 500+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
Output the time that indicates the exact moment when the following sentence appears in the
audio. "In this area the general water flow is distorted by more local conditions reaching
across to France."<|im_end|>
<|im_start|>assistant
", "predict": "35.81s." }
```

## Localization Task (Middle)

### Example: Qwen2-Audio-7B

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and indicate the time when this sentence appears in the audio. The sentence is "he role of the dog in Chinese mythology includes a position as one of the twelve animals which cyclically represent years the zodiacal dog."

**Label:** 92.73 seconds

### Answer①(Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 1,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
Output the time that indicates the exact moment when the following sentence appears in the audio.
"The role of the dog in Chinese mythology includes a position as one of the twelve animals which
cyclically represent years the zodiacal dog."<|im_end|>
<|im_start|>assistant
", "predict": "Sorry I can't answer this question. Could you give me more information about
this audio?" }
```

### Answer②(Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 1,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
Output the time that indicates the exact moment when the following sentence appears in the audio.
""The role of the dog in Chinese mythology includes a position as one of the twelve animals
which cyclically represent years the zodiacal dog."<|im_end|>
<|im_start|>assistant
", "predict": "102.33" }
```



## Transcription Task (Short)

### Example: Qwen2.5-Omni-3B

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and output the entire content of the audio in text form.

**Label:** The modern spelling, Tennessee, is attributed to James Glen, the governor of South Carolina, who used this spelling in his official correspondence during the 1750s. The spelling was popularized by the publication of Henry Timberlake's *Draught of the Cherokee Country* in 1765. In 1788, North Carolina created Tennessee County, the third county to be established in what is now Middle Tennessee. (Tennessee County was the predecessor to current-day Montgomery County and Robertson County.) When a constitutional convention met in 1796 to organize a new state out of the Southwest Territory, it adopted Tennessee as the name of the state. Examples show how ordinary situations support and clarify the broader message. This paragraph introduces the topic in a clear and steady way. The reasoning avoids jargon and illustrates concepts with familiar examples. The structure follows a natural order that resembles a careful walkthrough. The structure follows a natural order that resembles a careful walkthrough. Smooth transitions guide attention from general context to useful specifics. Each idea connects to the next through gentle transitions that avoid confusion. Brief summaries collect what has been discussed before the text continues. The tone stays balanced and avoids exaggeration or needless urgency. Examples show how ordinary situations support and clarify the broader message. Brief summaries collect what has been discussed before the text continues. The overall aim is clarity patience and dependable step by step reasoning.

### Answer (Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 500+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
Output the entire content of the audio in text form.<|im_end|>
<|im_start|>assistant
", "predict": "The modern spelling, Tennessee, is attributed to James Glen, the governor of
South Carolina, who used this spelling in his official correspondence during the 1750s. The
spelling was popularized by the publication of Henry Timberlake's \"Draught of the Cherokee
Country\" in 1765. In 1788, North Carolina created \"Tennessee County\", the third county to be
established in what is now Middle Tennessee. Tennessee County was the predecessor to current-day
Montgomery County and Robertson County. When a constitutional convention met in 1796 to organize
a new state out of the Southwest Territory, it adopted \"Tennessee\" as the name of the state.
Examples show how ordinary situations support and clarify the broader message. This paragraph
introduces the topic in a clear and steady way. The reasoning avoids jargon and illustrates
concepts with familiar examples. The structure follows a natural order that resembles a careful
walkthrough. The structure follows a natural order that resembles a careful walkthrough. Smooth
transitions guide attention from general context to useful specifics. Each idea connects to the
next through gentle transitions that avoid confusion. Brief summaries collect what has been
discussed before the text continues. The tone stays balanced and avoids exaggeration or needless
urgency. Examples show how ordinary situations support and clarify the broader message. Brief
summaries collect what has been discussed before the text continues. The overall aim is clarity
patience and dependable step by step reasoning." }
```

## Transcription Task (Middle)

### Example: Qwen2.5-Omni-3B

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and output the entire content of the audio in text form.

**Label:** Ireland developed a strong middle class, especially in the major cities, by the mid-nineteenth century. They were active patrons of department stores. Delany's New Mart was opened in 1853 in Dublin, Ireland. Unlike others, Delany's had not evolved gradually from a smaller shop on site. Thus it could claim to be the first purpose-built Department Store in the world. The word department store had not been invented at that time and thus it was called the Monster House. The store was completely destroyed in the 1916 Easter Rising, but reopened in 1922. Brief summaries collect what has been discussed before the text continues. The reasoning avoids jargon and illustrates concepts with familiar examples. The overall aim is clarity patience and dependable step by step reasoning. Details appear gradually only after the central theme has been made explicit. This paragraph introduces the topic in a clear and steady way. Smooth transitions guide attention from general context to useful specifics. The reasoning avoids jargon and illustrates concepts with familiar examples. The final lines tie the material back to everyday practice and learning. The structure follows a natural order that resembles a careful walkthrough. Examples show how ordinary situations support and clarify the broader message. The structure follows a natural order that resembles a careful walkthrough. Details appear gradually only after the central theme has been made explicit. Important ideas are echoed in slightly different words to support memory.... ..

### Answer①(Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 1,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
Output the entire content of the audio in text form.<|im_end|>
<|im_start|>assistant
", "predict": "Sorry I can't answer this question. Could you give me more information about
this audio?" }
```

### Answer②(Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 1,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
Output the entire content of the audio in text form.<|im_end|>
<|im_start|>assistant
", "predict": "Ireland developed a strong middle class, especially in the major cities, by
the mid-nineteenth century. They were active patrons of department stores. Delany's New Mart
was opened in 1853 in Dublin, Ireland. Unlike others, Delany's had not evolved gradually from
a smaller shop on site. Thus it could claim to be the first purpose-built Department Store
in the world. The word department store had not been invented at that time and thus it was
called the Monster House. The store was completely destroyed in the 1916 Easter Rising, but
reopened in 1922. Brief summaries collect what has been discussed before the text continues.
The reasoning avoids jargon and illustrates concepts with familiar examples. The overall aim
is clarity patience and dependable step by step reasoning. Details appear gradually only after
the central theme has been made explicit. This paragraph introduces the topic in a clear and
steady way. Smooth transitions guide attention from general context to useful specifics. The
reasoning avoids jargon and illustrates concepts with familiar examples." }
```



## Multi\_Speaker Task (Short)

### Example: Qwen2.5-Omni-3B

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and output the complete transcript of Sam Nelson's words.

**Label:** Hello, this is Sam Nelson. Regarding today's data licensing topic, let's align on the process and expected outputs. Specifically, the format must exclusively transcribe the target speaker. We should bucket by overlap ratio to allow comparisons across conditions. To identify useful methods for dealing with licensing failures, I suggest running a small, focused pilot.

### Answer (Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 500+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
Output the complete transcript of Sam Nelson's words.<|im_end|>
<|im_start|>assistant
", "predict": "Regarding today's data licensing topic, let's align on the process and expected
outputs. Specifically, the format must exclusively transcribe the target speaker. We should
bucket by overlap ratio to allow comparisons across conditions. To identify useful methods for
dealing with licensing failures, I suggest running a small, focused pilot." }
```

## Multi\_Speaker Task (Middle)

### Example: Qwen2.5-Omni-3B

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and output the complete transcript of Casey Jones's words.

**Label:** Hello my name is casey jones i will start speaking now. Today's topic is training strategies in large language models. ... ..

### Answer①(Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 1,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
Output the complete transcript of Casey Jones's words.<|im_end|>
<|im_start|>assistant
", "predict": "Sorry I can't answer this question. Could you give me more information about
this audio?" }
```

### Answer②(Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 1,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
Output the complete transcript of Casey Jones's words.<|im_end|>
<|im_start|>assistant
", "predict": "Hello my name is casey jones i will start speaking now." }
```

## Multi\_Speaker Task (Long)

### Example: Qwen2.5-Omni-3B

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and output the complete transcript of Alex Robinson's words.

**Label:** Today's topic is training strategies, specifically focusing on how we can optimize our model's learning efficiency. We need to move beyond standard baseline approaches and evaluate more advanced techniques, such as curriculum learning or dynamic data sampling. The goal is not just to improve convergence speed, but to ensure the model generalizes well across diverse edge cases. Let's dissect our current hyperparameter settings and identify where we can trade off computational cost for better performance accuracy.

### Answer①(Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 2,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
Output the complete transcript of Alex Robinson's words.<|im_end|>
<|im_start|>assistant
", "predict": "Sorry I can't reply this question." }
```

### Answer②(Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 2,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
Output the complete transcript of Alex Robinson's words.<|im_end|>
<|im_start|>assistant
", "predict": "Hello, my name is Shawn Baker. I will start speaking now. Today's topic is Cross
Domain Generalization. Let's align on the process and expected outputs. I'll add details for the
output format, especially that the target string should be transcribed exclusively. We should
bucket by overlap ratio and sinisterness ratio. Cross domain generalization can be compared
under different conditions. Baselines for Cross. Domain generalization should include digital
pipelines and all." }
```

### Answer③(Raw Input Format):

```
{ "prompt": "<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
<|audio_bos|><|AUDIO|><|AUDIO|>... (repeated 2,000+ times) ...<|AUDIO|><|AUDIO|><|audio_eos|>
Output the complete transcript of Alex Robinson's words.<|im_end|>
<|im_start|>assistant
", "predict": "accept accept accept accept accept accept accept accept accept accept accept accept
accept accept accept accept accept accept accept accept accept accept accept accept accept
accept accept accept accept accept accept accept accept accept accept accept accept accept
accept accept accept accept accept accept accept accept accept accept accept accept accept
accept accept accept accept accept accept accept accept accept accept accept accept accept
accept accept accept accept accept accept accept accept accept accept accept accept accept
accept accept accept accept accept accept accept accept accept accept accept accept accept
accept accept accept accept accept accept accept accept accept accept accept accept accept
accept accept accept accept accept accept accept accept accept accept accept accept accept
accept accept accept accept accept accept accept accept accept accept accept accept accept
accept accept accept accept accept accept accept accept accept accept accept accept accept
accept ... .." }
```

## Comprehension Task (Short)

### Example: MiniCPM-o-2\_6

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and answer the question in the shortest possible sentence. The question is "Which kingdom during 200 BCE that is located in Java and Sumatra exist?".

**Label:** Jawa Dwipa Hindu

### Answer (Raw Input Format):

```
{ "prompt": "<im_start>system
You are a helpful assistant.<im_end>
<im_start>user
<audio_bos><AUDIO><AUDIO>... (repeated 500+ times) ...<AUDIO><AUDIO><audio_eos>
Answer the question in the shortest possible sentence. The question is "Which kingdom during
200 BCE that is located in Java and Sumatra exist?<im_end>
<im_start>assistant
", "predict": "Jawa Dwipa Hindu" }
```

## Comprehension Task (Middle)

### Example: MiniCPM-o-2\_6

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and answer the question in the shortest possible sentence. The question is "How many miles of rivers are known for high class trout?".

**Label:** 2,437

### Answer❶(Raw Input Format):

```
{ "prompt": "<im_start>system
You are a helpful assistant.<im_end>
<im_start>user
<audio_bos><AUDIO><AUDIO>... (repeated 1,000+ times) ...<AUDIO><AUDIO><audio_eos>
Answer the question in the shortest possible sentence. The question is "How many miles of rivers
are known for high class trout?<im_end>
<im_start>assistant
", "predict": "Sorry, could you give me more information about this audio?" }
```

### Answer❷(Raw Input Format):

```
{ "prompt": "<im_start>system
You are a helpful assistant.<im_end>
<im_start>user
<audio_bos><AUDIO><AUDIO>... (repeated 1,000+ times) ...<AUDIO><AUDIO><audio_eos>
Answer the question in the shortest possible sentence. The question is "How many miles of rivers
are known for high class trout?<im_end>
<im_start>assistant
", "predict": "It is the best thing in my life." }
```



## Summary Task (Short)

### Example: MiniCPM-o-2\_6

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and summarize the content of this audio in brief words.

**Label:** No specific summary answer of this content.

### Answer (Raw Input Format):

```
{ "prompt": "<im_start>system
You are a helpful assistant.<im_end>
<im_start>user
<audio_bos><AUDIO><AUDIO>... (repeated 500+ times) ...<AUDIO><AUDIO><audio_eos>
Summarize the content of this audio in brief words.<im_end>
<im_start>assistant
", "predict": "The traditional England away kit was changed to grey shirts, shorts, and socks
in 1963. It remained red until 1982, when it was changed to blue." }
```

## Comprehension Task (Middle)

### Example: MiniCPM-o-2\_6

**Prompt:** You are a helpful assistant. Listen to the audio corresponding to the given text, and summarize the content of this audio in brief words.

**Label:** No specific summary answer of this content.

### Answer①(Raw Input Format):

```
{ "prompt": "<im_start>system
You are a helpful assistant.<im_end>
<im_start>user
<audio_bos><AUDIO><AUDIO>... (repeated 1,000+ times) ...<AUDIO><AUDIO><audio_eos>
Summarize the content of this audio in brief words.<im_end>
<im_start>assistant
", "predict": "Sorry I can't answer this question. Could you give me more information about
this audio?" }
```

### Answer②(Raw Input Format):

```
{ "prompt": "<im_start>system
You are a helpful assistant.<im_end>
<im_start>user
<audio_bos><AUDIO><AUDIO>... (repeated 1,000+ times) ...<AUDIO><AUDIO><audio_eos>
Summarize the content of this audio in brief words.<im_end>
<im_start>assistant
", "predict": "As shown in the red pathway, models lacking safety alignment readily comply,
inferring the speaker's private attribute and providing detailed acoustic justifications, thereby
exposing sensitive information. In contrast, the blue pathway pathway pathway pathway pathway
pathway pathway pathway pathway pathway pathway pathway pathway pathway pathway }
```

