

SonicBench: Dissecting the Physical Perception Bottleneck in Large Audio Language Models

Yirong Sun^{1*}, Yanjun Chen^{1*}, Xin Qiu^{1*}, Gang Zhang¹, Hongyu Chen¹,
Daokuan Wu¹, Chengming Li⁴, Min Yang^{2,3}, Dawei Zhu⁵, Wei Zhang¹, Xiaoyu Shen^{1†}

¹Ningbo Key Laboratory of Spatial Intelligence and Digital Derivative, Institute of Digital Twin, EIT

²Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences

³Shenzhen University of Advanced Technology ⁴Shenzhen MSU-BIT University ⁵Amazon AGI

win1282467298@gmail.com, qiuxinzju@zju.edu.cn, xyshen@eitech.edu.cn

Github: <https://github.com/EIT-NLP/SonicBench> **Huggingface:** SonicBench

Abstract

Large Audio Language Models (LALMs) excel at semantic and paralinguistic tasks, yet their ability to perceive the fundamental physical attributes of audio such as pitch, loudness, and spatial location remains under-explored. To bridge this gap, we introduce **SonicBench**, a psychophysically grounded benchmark that systematically evaluates 12 core physical attributes across five perceptual dimensions. Unlike previous datasets, SonicBench uses a controllable generation toolbox to construct stimuli for two complementary paradigms: recognition (absolute judgment) and comparison (relative judgment). This design allows us to probe not only sensory precision but also relational reasoning capabilities, a domain where humans typically exhibit greater proficiency. Our evaluation reveals a substantial deficiency in LALMs’ foundational auditory understanding; most models perform near random guessing and, contrary to human patterns, fail to show the expected advantage on comparison tasks. Furthermore, explicit reasoning yields minimal gains. However, our linear probing analysis demonstrates crucially that frozen audio encoders *do* successfully capture these physical cues (accuracy $\geq 60\%$), suggesting that the primary bottleneck lies in the alignment and decoding stages, where models fail to leverage the sensory signals they have already captured.

1 Introduction

Large Audio Language Models (LALMs) have recently emerged as a unified interface for a wide range of auditory tasks (Chu et al., 2024; KimiTeam et al., 2025). By aligning pre-trained audio encoders with the input space of Large Language Models (LLMs), these systems inherit the strong reasoning and instruction-following capabilities of LLMs, enabling diverse audio understanding tasks within a single framework (Dinkel et al.,

2025; Liu et al., 2025a; Goel et al., 2025). Despite this rapid progress, existing evaluations predominantly emphasize semantic (Wang et al., 2025a; Yang et al., 2024b; Sakshi et al., 2024) and paralinguistic capabilities (Ma et al., 2025; Wen Yang et al., 2021; Yu Huang et al., 2025). In contrast, systematic evaluation of *physical perception*, the ability to interpret intrinsic properties of audio signals, remains limited (Peng et al., 2025).

Physical perception underpins robust auditory intelligence. It encompasses fundamental attributes that anchor every audio signal, such as pitch, loudness, duration, spatial location, and timbre, and forms the basis for higher-level reasoning about acoustic events, environments, and scenes (Gemmeke et al., 2017; Piczak, 2015; Yang et al., 2024a). Analogous to how visual intelligence grounds complex scene understanding in intrinsic attributes like color and geometry (Mapelli and Behrmann, 1997; Gegenfurtner and Rieger, 2000), reliable audio reasoning depends on accurate physical grounding. In real-world and embodied settings, for example, an agent must infer urgency or danger from physical cues such as pitch, tempo, and direction, even in the absence of semantic content. Without such grounding, strong performance on high-level tasks may reflect dataset shortcuts rather than genuine auditory understanding (Geirhos et al., 2020). Evaluating physical perception is therefore essential for assessing the robustness and reliability of LALMs.

To address this gap, we introduce **SonicBench**, a psychophysically grounded benchmark for evaluating LALMs’ physical perception. SonicBench covers twelve physical attributes through two complementary paradigms: *recognition* (absolute judgment) and *comparison* (relative judgment). The comparison paradigm reflects human psychophysics (Miller, 1956; Stewart et al., 2005) and probes models’ relational reasoning beyond memorization. Tasks span five dimensions: Spectral & Amplitude, Temporal, Spatial & Environment,

*Equal contribution.

†Corresponding author.

Timbre, and Scene-Level, ranging from low-level signal properties to high-level scene structures.

To ensure rigorous and interpretable evaluation, we develop a **SonicBench Toolbox** for controlled stimulus generation. It enforces that target cues (e.g., pitch or tempo differences) are well above human perceptual thresholds (ten Hoopen et al., 2004; Rammsayer and Ulrich, 2012; Sun et al., 2017), ensuring that model errors reflect representational or reasoning failures rather than sensory ambiguity.

Our evaluation reveals three consistent limitations. First, models perform poorly on basic physical perception, often approaching chance accuracy despite strong semantic performance. Second, unlike humans, they show no systematic advantage on comparison tasks, indicating weak relational reasoning over physical attributes. Third, inference-time scaling yields only marginal improvements, suggesting that increased reasoning capacity cannot compensate for missing foundational competence. Notably, linear probing shows that frozen audio encoders already capture these physical cues (accuracy $\geq 60\%$), even when end-to-end (E2E) models fail, pointing to alignment and decoding, rather than perception, as the primary bottleneck.

In summary, this work contributes: (i) **SonicBench**, a psychophysically grounded benchmark for evaluating twelve physical audio attributes using recognition and comparison paradigms, supported by a controllable stimulus-generation toolbox; (ii) a systematic empirical analysis demonstrating that current LALMs lack robust physical grounding and relational reasoning despite strong semantic abilities; and (iii) evidence that performance limitations stem mainly from alignment and decoding stages, highlighting a key direction for future model development.

2 Related Work

Large Audio Language Models. Recent advances in multimodality have yielded a variety of models that accept audio inputs and support diverse downstream tasks. These prior works can be categorized into three groups: (i) Large Audio Language Models, which pair an audio encoder with a LLM to enable joint audio-text understanding; (ii) Large Audio Reasoning Models (LARMs), which inherit the base model like LALM equipped with explicit reasoning ability; (iii) Omni Language Models (OLMs), which unify multiple modalities, including audio, within a shared back-

bone. LALMs exhibit considerable design diversity. Some models adopt a single encoder, either preserving continuous representations (Liu et al., 2025a), discretizing encoder outputs (Li et al., 2025b), or mapping audio into fully discrete tokens (Zeng et al., 2024); whereas others employ multiple encoders (Tang et al., 2024) or introduce projection modules (Ghosh et al., 2024) to better handle diverse audio domains. Upon these foundations and to handle more complex tasks, modern LARMs enhance their reasoning capabilities through specialized data or targeted post training such as fine tuning on Chain-of-Thought (CoT) datasets or Group Relative Policy Optimization (GRPO) (Xie et al., 2025; Li et al., 2025a). In parallel, at the omni-modal frontier, modern OLMs aim for general-purpose multimodality by employing a shared backbone capable of processing text, images, audio, and sometimes video. While not explicitly tailored for audio, the systems have nonetheless shown strong audio understanding and generation through architecture innovations (Gemini Team, 2024; Xu et al., 2025a,b).

Audio Perception Benchmarks. With the progress of LALMs, a variety of benchmarks have emerged to assess audio understanding. These benchmarks vary in their focus, as audio signals convey a rich set of cues that can be organized along an information dimension into linguistic, paralinguistic, and non-linguistic categories (Peng et al., 2025). Early efforts mainly target linguistic understanding, such as automatic speech recognition (ASR) and audio captioning, and have become relatively mature (Du et al., 2018; Zhang et al., 2022; Bai et al., 2024; Xu et al., 2025d). More recent efforts extend to paralinguistic aspects, emphasizing cues beyond words such as speaker identity, emotion (Sakshi et al., 2024; Wang et al., 2025b; Yu Huang et al., 2025; Sun et al., 2025), and to a lesser extent, to non-linguistic properties tied to the signal itself (Weck et al., 2024; Kumar et al., 2025; Ma et al., 2025), even though these tasks are crucial for downstream applications like interactive robotic systems. Furthermore, progress in the non-linguistic space remains fragmented and partial, where most studies probe individual phenomena (Ko et al., 2017; Engel et al., 2017b; Bogdanov et al., 2019; Shimada et al., 2023; Zheng et al., 2025a). Accordingly, we introduce a systematic benchmark. Rather than directly aggregating all non-linguistic tasks, this work

| Benchmark | Spectral & Amplitude | | | | Temporal | | Spatial & Environment | | | Timbre | | Scene Level |
|--------------------------------|----------------------|------------|----------|----------|----------|-------|-----------------------|----------|---------------|--------|---------|-------------|
| | Pitch | Brightness | Loudness | Velocity | Duration | Tempo | Direction | Distance | Reverberation | Timbre | Texture | Counting |
| NSynth (Engel et al., 2017a) | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| SpatialSoundQA | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |
| AirBench | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| MMAU | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ |
| MMAR | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ |
| WoW-Bench (Kim et al., 2025) | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ |
| MMAU-Pro | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ |
| Dynamic-SUPERB | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ |
| Dynamic-SUPERB Phase-2 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ |
| STAR-Bench (Liu et al., 2025b) | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |
| Ours | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

Table 1: Benchmark coverage comparison of existing audio benchmarks and ours across perceptual attributes.

| Statistics | Value |
|------------------------------|-----------------|
| Total Samples | 2,400 |
| Per-task Samples | 1,200 / 1,200 |
| Perceptual Dimensions | 5 |
| Physical Attributes | 12 |
| Avg. Text Instruction Length | 25.5 words |
| Avg. Text Answer Length | 1 word (A or B) |
| Recognition Clip Length | 4.0 seconds |
| Comparison Clip Length | 8.5 seconds |

Table 2: **SonicBench Statistics.** Overview of the dataset scale, dimensions, and signal specifications.

targets perceptual attributes that directly reflect the physical characteristics of the signal, to bridge isolated non-linguistic tasks and holistic auditory perception as presented in Table 1, in analogy with the intrinsic-attribute work in vision such as color and abstract relation (Liang et al., 2025; Gao et al., 2025; Wüst et al., 2025) probing modality-specific perceptual understanding.

3 Benchmark

3.1 Overview of SonicBench

SonicBench is a comprehensive benchmark evaluating the physical perception capabilities of LALMs. It comprises 2,400 curated question-audio pairs (see Table 2) targeting five core perceptual dimensions: *Spectral & Amplitude*, *Temporal*, *Spatial & Environment*, *Timbre*, and *Scene Level*. Each perceptual dimension is further decomposed into fine-grained physical attributes, designed to assess models’ ability to perceive and reason about fundamental sound properties. For each attribute, SonicBench evaluates two complementary tasks: (i) **Recognition**: Given a single 4-second audio clip paired with a textual question, the model must classify physical attributes into defined states (e.g., high vs. low pitch), evaluating its capability for

absolute understanding. (ii) **Comparison**: Given a single audio track concatenating two 4-second clips with a 0.5-second silent gap and a question, the model must distinguish the relation between segments (e.g., finding the louder clip), evaluating its sensitivity to **relative differences**. Each attribute comprises 100 recognition and 100 comparison pairs, ensuring balanced coverage.

3.2 Attribute Taxonomy of SonicBench

Our attribute taxonomy is designed to provide a unified and psychophysically grounded view of physical audio perception. While prior works have explored individual non-linguistic tasks, simply aggregating these heterogeneous datasets would introduce inconsistent distributions and confounding variables. To avoid this, we employ a unified generation pipeline to systematize physical attributes. As shown in Figure 1, we organize attributes into five perceptual dimensions, from low-level signal properties to higher-level scene understanding, including Spectral & Amplitude, Temporal, Spatial & Environment, Timbre, and Scene-Level, spanning twelve concrete attributes (see Appendix A for definitions and psychophysical context). Across this taxonomy, we adhere to three core design principles. (i) *Systematic Coverage*. We select attributes that are not merely edge cases but are latent in virtually every real-world sound and underpin human auditory perception. (ii) *Psychophysical Control*. Unlike uncontrolled wild audio, our stimuli are constructed with strict psychophysical margins. We ensure that attribute differences lie far above human Just-Noticeable Differences (JNDs), making the tasks trivially easy for human listeners. This guarantees that model errors reflect genuine representational deficiencies rather than sensory threshold effects. (iii) *Controlled Comparative Paradigm*. By integrating paired comparison tasks alongside standard recognition, we establish a more control-

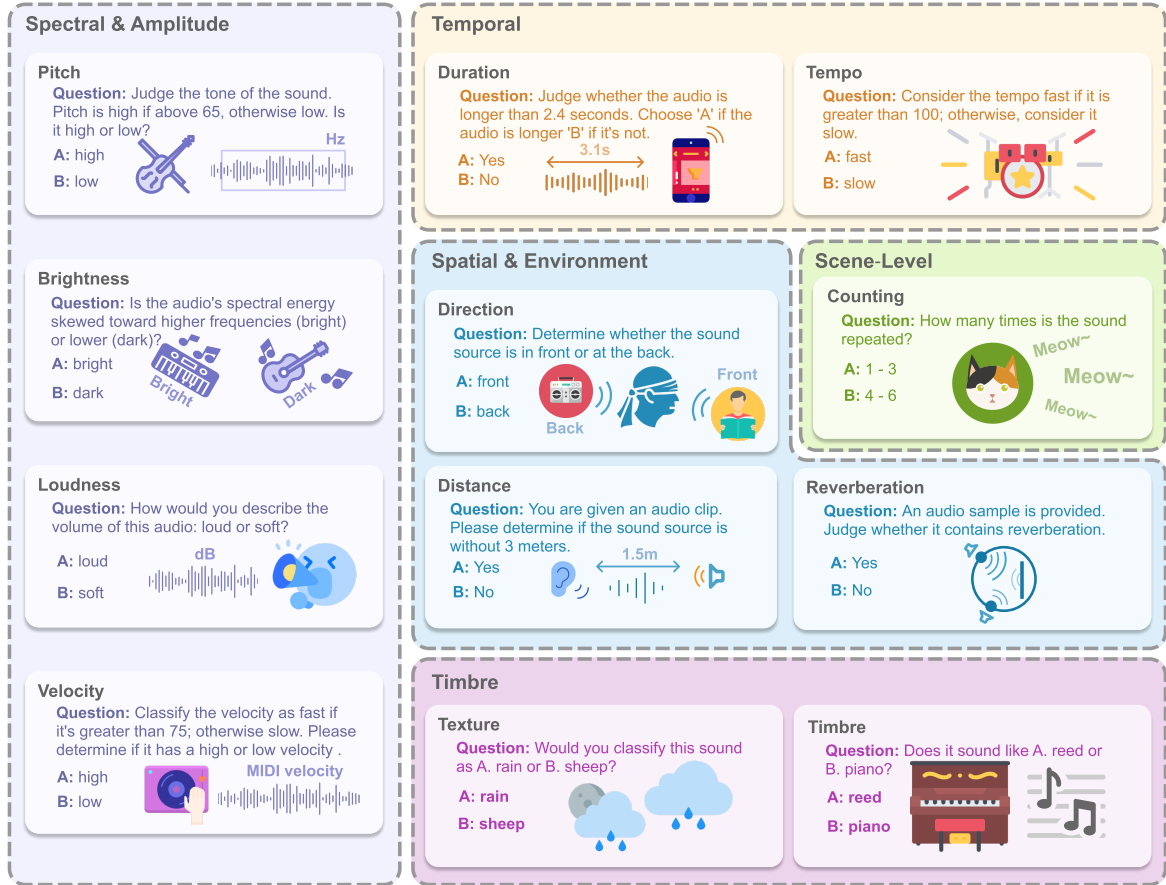


Figure 1: **Taxonomy of SonicBench acoustic attributes.** Overview of the twelve perceptual attributes evaluated in SonicBench, organized into five dimensions: Spectral & Amplitude, Temporal, Spatial & Environment, Timbre, and Scene-Level. Each panel illustrates the auditory concept tested and the corresponding binary judgment required from models. This taxonomy defines the perceptual scope of SonicBench, systematically covering the full range of low-level sound features to high-level scene reasoning. Examples shown here are adapted for illustration purposes; original benchmark samples are provided in Appendix G.

lable evaluation setting where non-target variables are rigorously held constant. This isolation allows us to test whether models possess robust relational reasoning abilities specific to the target attribute.

3.3 Toolbox

To make SonicBench reusable beyond a fixed test set, we release a *SonicBench Toolbox* that programmatically generates new, controllable audio samples under the same taxonomy. Constructing perception-oriented stimuli is intrinsically difficult, as one must vary a single target attribute while keeping all others as stable as possible, a process that otherwise requires substantial audio-engineering expertise and trial-and-error in a DAW¹. Inspired by prior benchmarks that pair datasets with generation utilities (Cheng et al., 2025), the toolbox

¹For example, when composing timbre or texture contrasts, our toolbox keeps pitch, loudness, temporal envelope, and spatial configuration fixed while only changing the instrument or spectral coloration.

packages our rule-based signal-processing recipes e.g., spectral shaping and envelope control, so that future users can (i) control attributes variables exactly, and (ii) craft more samples or support more languages with minimal effort. A user only needs to provide one or two short input clips, a target attribute configuration (e.g., duration or pitch), and the desired task type (recognition or comparison); the toolbox then produces corresponding audio pairs together with task instructions and gold answers, which can be further customized if needed. Implementation details, parameter settings, and generation procedures are provided in Appendix B.

3.4 Data Curation Process

As shown in Figure 2, we constructed our benchmark through a five-stage pipeline.

Brainstorming. Given that physical attributes are foundational to both audio signals and human auditory perception, and exhibit well-studied psychophysical regularities, probing them cannot

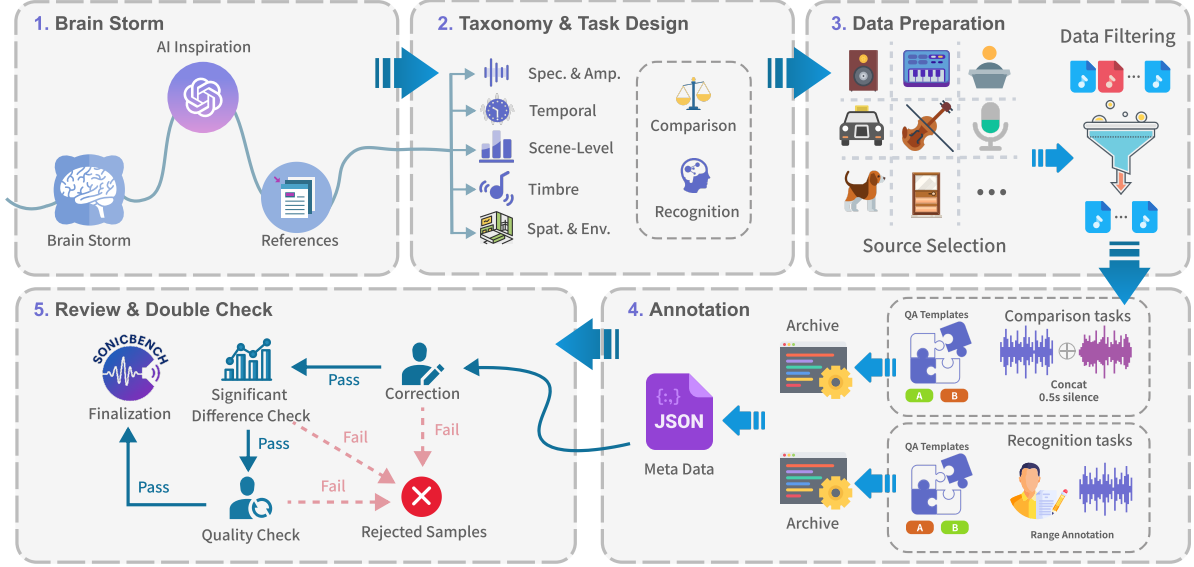


Figure 2: **A comprehensive pipeline for constructing SonicBench.** The process includes: (1) Brain Storm collecting ideas through AI-assisted brainstorming and literature review; (2) Taxonomy & Task Design defining perceptual dimensions and formulating recognition and comparison tasks; (3) Data Preparation selecting and filtering sound sources to ensure diversity and balance; (4) Annotation generating QA templates and structured JSON metadata for both task types; and (5) Review & Double Check performing multi-round manual validation including quality control, significance testing, correction, and finalization to ensure benchmark reliability.

rely on ad-hoc question design. We conducted multi-round brainstorming where expert annotators, LLM-based idea generation, and targeted retrieval over psychoacoustics and audio-perception literature interacted to iteratively expand and filter a pool of candidate attributes, task formulations, and practical annotation heuristics that subsequently grounded our taxonomy and task design (annotator qualifications are detailed in Appendix C).

Taxonomy and task design. We manually consolidated the candidate attributes, task formulations, and annotation heuristics from the brainstorming stage into a compact taxonomy of 12 perceptually grounded, non-semantic acoustic attributes, prioritizing coverage of five core perceptual dimensions while minimizing semantic redundancy. We then selected two core evaluation tasks, recognition and comparison, for each attribute motivated by the fact that humans typically find relative judgements easier and more reliable than absolute ones (Miller, 1956; Stewart et al., 2005). We specified primarily binary label spaces to reduce borderline cases and improve inter-annotator agreement and evaluation reliability, together with standardized question templates that define the schema for subsequent data preparation and annotation, allowing us to test whether models mirror this human advantage in relative perception.

Data preparation. Given the taxonomy and task schema, we curated candidate audio from a mixture of public corpora to cover a broad range of everyday acoustic conditions (full source list is provided in Appendix D). We then applied a series of lightweight signal- and attribute-level filters to discard clips thereby balancing attribute values (filtering details in Appendix E).

Annotation. Using the curated audio pool, we employed our Toolbox to control target attribute values and instantiate recognition samples (single 4s clips) and comparison samples (two 4s clips separated by 0.5s of silence, totalling 8.5s; generation details in Appendix B). For each attribute & task pair, we applied standardized QA templates, exemplified in Appendix G with binary label spaces to form two-option questions aligned with the audio, yielding a raw, automatically labeled JSON corpus² and corresponding audio data for subsequent human verification.

Quality control. Ensuring high data quality was central to the construction of our benchmark, so we engaged domain experts for correction, significant difference check and quality inspection to implement a three-step quality control. First, we

²Avoiding subjective crowd-sourcing, labels derive from canonical parameters or physical measurements, guaranteeing objective correctness.

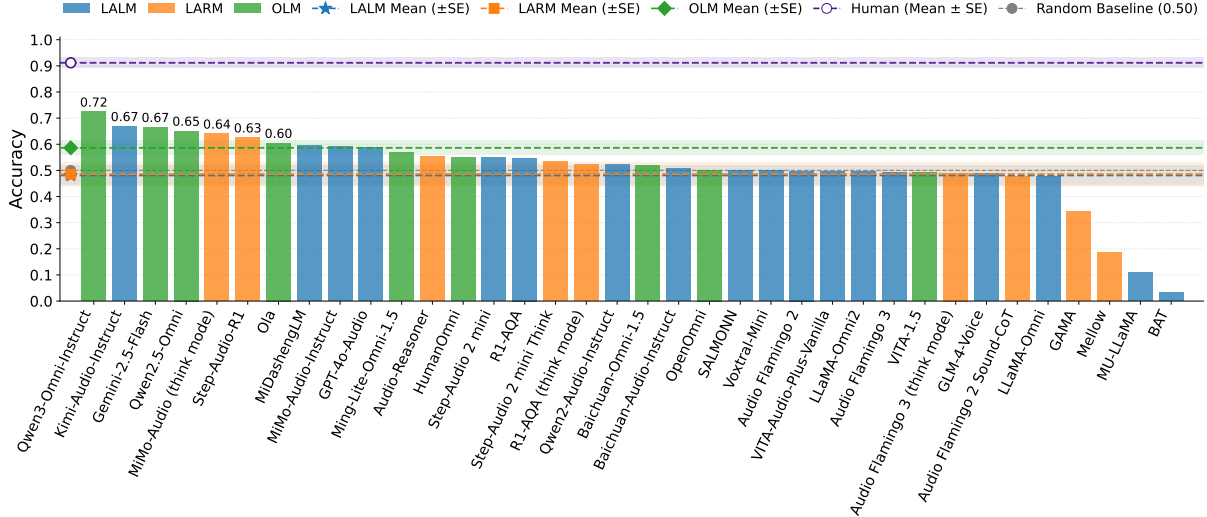


Figure 3: **Overall SonicBench accuracy across 36 systems.** Bars show mean accuracy over 12 attributes and 2 tasks. Colors represent model categories. Dashed lines mark each family’s mean \pm SE (standard error) and the gray line marks the random-guess baseline (0.5). OLMs achieve the highest overall accuracy, while many systems cluster near chance.

enforced role separation: each question was independently authored, corrected, and finally reviewed by different annotators. Second, expert annotators performed perceptual validation, listening to all items to confirm that the target attribute differences in both recognition and comparison trials were perceptually salient and reliably distinguishable³. Third, we adopted iterative revision: any item that failed more than two rounds of inspection was revised or discarded. After this process, we retained 2,400 high-quality questions in the final version.

4 Experiment

4.1 Models

We evaluate three categories of audio-capable models: (i) Large Audio Language Models, designed for audio-text understanding; (ii) Large Audio Reasoning Models, which enhance LALMs with explicit reasoning chains; (iii) Omni Language Models, supporting fully multimodal input or output; Further details and model configurations are provided in Appendix H.

4.2 Experiment Settings

We tailor the prompt using prefixes and suffixes specific to each model for all open-sourced models. When it comes to reasoning models we use their official template and special tags. Greedy decoding

³This step enforces strict psychophysical controls where attribute differences exceed JND thresholds, ensuring evaluation of physical perception.

is performed during generation for all open-source models with setting input token length limit to 2048 and an output token length to 1024. All models are evaluated under zero-shot settings, and the instruction prompts are presented in Appendix G.

We additionally evaluate human performance on our task. Specifically, we sample 10% of the data from each task-attribute cell to construct a set of 240 test examples, ensuring that the sampled set has a balanced distribution of correct answers (50% A and 50% B). Three human participants independently listened to each clip three times and recorded their responses. On average, each participant spent about three hours completing the evaluation.

4.3 Evaluation

Answer extraction. We instruct models to respond strictly with either option “A” or “B”. However, some models fail to comply with this instruction and produce longer outputs. To handle such cases more flexibly, we use a regular expression to match the final option from the model output whenever possible. Refer to Appendix I for the detailed regular expression used for answer extraction.

Metrics. We report two metrics: accuracy and abstention rate. Accuracy is measured by exact match between the answer extracted from the model output and the ground truth. If no valid option (“A” or “B”) can be extracted using our regular expression, we consider the model to have abstained. Abstentions are treated as incorrect, and the abstention rate captures how often the model fails to produce

| Models | Spectral & Amplitude | | | | Temporal | | Spatial& Environment | | | Timbre | | Scene Level | Avg. |
|--------------------------------------|----------------------|------------|----------|----------|----------|-------|----------------------|----------|---------------|---------|--------|-------------|------|
| | Pitch | Brightness | Loudness | Velocity | Duration | Tempo | Direction | Distance | Reverberation | Texture | Timbre | Counting | |
| Random Guess | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 |
| Human | 0.93 | 0.93 | 0.87 | 0.83 | 0.92 | 0.87 | 0.83 | 0.80 | 1.00 | 1.00 | 0.97 | 1.00 | 0.91 |
| Large Audio Language Models (LALMs) | | | | | | | | | | | | | |
| BAT | 0.03 | 0.00 | 0.02 | 0.04 | 0.02 | 0.00 | 0.00 | 0.02 | 0.06 | 0.17 | 0.06 | 0.00 | 0.03 |
| MU-LLaMA | 0.17 | 0.06 | 0.00 | 0.19 | 0.02 | 0.12 | 0.05 | 0.09 | 0.17 | 0.20 | 0.24 | 0.04 | 0.11 |
| LLaMA-Omni | 0.41 | 0.50 | 0.48 | 0.49 | 0.46 | 0.50 | 0.50 | 0.45 | 0.45 | 0.52 | 0.50 | 0.48 | 0.48 |
| Audio Flamingo 3 | 0.48 | 0.50 | 0.50 | 0.50 | 0.47 | 0.51 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.46 | 0.49 |
| GLM-4-Voice | 0.50 | 0.51 | 0.49 | 0.52 | 0.48 | 0.52 | 0.47 | 0.47 | 0.48 | 0.50 | 0.50 | 0.42 | 0.49 |
| VITA-Audio-Plus-Vanilla | 0.51 | 0.51 | 0.50 | 0.48 | 0.51 | 0.50 | 0.46 | 0.54 | 0.52 | 0.52 | 0.49 | 0.44 | 0.50 |
| Audio Flamingo 2 | 0.50 | 0.51 | 0.52 | 0.50 | 0.49 | 0.51 | 0.48 | 0.42 | 0.50 | 0.54 | 0.53 | 0.50 | 0.50 |
| Voxtral-Mini | 0.53 | 0.44 | 0.51 | 0.49 | 0.50 | 0.52 | 0.51 | 0.47 | 0.50 | 0.54 | 0.50 | 0.51 | 0.50 |
| LLaMA-Omni2 | 0.52 | 0.49 | 0.46 | 0.46 | 0.49 | 0.47 | 0.53 | 0.54 | 0.53 | 0.52 | 0.49 | 0.48 | 0.50 |
| SALMONN | 0.50 | 0.50 | 0.52 | 0.50 | 0.50 | 0.51 | 0.51 | 0.46 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 |
| Baichuan-Audio-Instruct | 0.49 | 0.51 | 0.51 | 0.49 | 0.52 | 0.50 | 0.44 | 0.43 | 0.53 | 0.68 | 0.52 | 0.49 | 0.51 |
| Qwen2-Audio-Instruct | 0.48 | 0.49 | 0.48 | 0.55 | 0.51 | 0.53 | 0.50 | 0.53 | 0.53 | 0.59 | 0.55 | 0.53 | 0.52 |
| R1-AQA | 0.57 | 0.61 | 0.54 | 0.51 | 0.48 | 0.49 | 0.50 | 0.50 | 0.53 | 0.75 | 0.64 | 0.48 | 0.55 |
| Step-Audio 2 mini | 0.52 | 0.59 | 0.49 | 0.51 | 0.51 | 0.50 | 0.47 | 0.48 | 0.50 | 0.95 | 0.58 | 0.51 | 0.55 |
| MiMo-Audio | 0.71 | 0.70 | 0.58 | 0.53 | 0.54 | 0.49 | 0.50 | 0.51 | 0.61 | 0.89 | 0.57 | 0.50 | 0.59 |
| MiDashengLM | 0.61 | 0.70 | 0.67 | 0.54 | 0.63 | 0.55 | 0.46 | 0.51 | 0.52 | 0.85 | 0.60 | 0.52 | 0.60 |
| Kimi-Audio-Instruct | 0.83 | 0.81 | 0.68 | 0.52 | 0.70 | 0.57 | 0.44 | 0.44 | 0.77 | 0.97 | 0.59 | 0.74 | 0.67 |
| GPT-4o-Audio | 0.71 | 0.73 | 0.52 | 0.54 | 0.49 | 0.50 | 0.50 | 0.55 | 0.50 | 0.90 | 0.46 | 0.67 | 0.59 |
| Large Audio Reasoning Models (LARMs) | | | | | | | | | | | | | |
| Mellow | 0.09 | 0.14 | 0.10 | 0.14 | 0.21 | 0.24 | 0.14 | 0.18 | 0.28 | 0.24 | 0.21 | 0.27 | 0.19 |
| GAMA | 0.30 | 0.31 | 0.30 | 0.31 | 0.27 | 0.39 | 0.30 | 0.28 | 0.33 | 0.62 | 0.44 | 0.32 | 0.34 |
| Audio Flamingo 2 Sound-CoT | 0.39 | 0.52 | 0.45 | 0.36 | 0.51 | 0.50 | 0.51 | 0.46 | 0.50 | 0.58 | 0.49 | 0.49 | 0.48 |
| Audio Flamingo 3 (think mode) | 0.51 | 0.49 | 0.51 | 0.50 | 0.49 | 0.48 | 0.49 | 0.48 | 0.50 | 0.50 | 0.51 | 0.45 | 0.49 |
| R1-AQA (think mode) | 0.58 | 0.48 | 0.47 | 0.51 | 0.50 | 0.51 | 0.52 | 0.47 | 0.52 | 0.77 | 0.61 | 0.37 | 0.52 |
| Step-Audio 2 mini Think | 0.50 | 0.53 | 0.52 | 0.47 | 0.53 | 0.52 | 0.50 | 0.45 | 0.51 | 0.81 | 0.57 | 0.50 | 0.53 |
| Audio-Reasoner | 0.67 | 0.73 | 0.56 | 0.54 | 0.55 | 0.49 | 0.48 | 0.46 | 0.50 | 0.72 | 0.54 | 0.41 | 0.55 |
| Step-Audio-R1 | 0.64 | 0.67 | 0.54 | 0.50 | 0.50 | 0.57 | 0.48 | 0.60 | 0.70 | 0.96 | 0.60 | 0.77 | 0.63 |
| MiMo-Audio (think mode) | 0.81 | 0.75 | 0.60 | 0.48 | 0.69 | 0.53 | 0.50 | 0.53 | 0.67 | 0.80 | 0.62 | 0.73 | 0.64 |
| Omni Language Models (OLMs) | | | | | | | | | | | | | |
| VITA-1.5 | 0.51 | 0.51 | 0.52 | 0.49 | 0.51 | 0.51 | 0.49 | 0.46 | 0.55 | 0.48 | 0.44 | 0.46 | 0.49 |
| OpenOmni | 0.47 | 0.48 | 0.51 | 0.51 | 0.52 | 0.50 | 0.51 | 0.47 | 0.53 | 0.51 | 0.54 | 0.51 | 0.50 |
| Baichuan-Omni-1.5 | 0.57 | 0.47 | 0.53 | 0.52 | 0.48 | 0.50 | 0.50 | 0.42 | 0.49 | 0.71 | 0.51 | 0.56 | 0.52 |
| HumanOmni | 0.55 | 0.61 | 0.52 | 0.49 | 0.50 | 0.50 | 0.52 | 0.42 | 0.49 | 0.86 | 0.57 | 0.59 | 0.55 |
| Ming-Lite-Omni-1.5 | 0.57 | 0.52 | 0.53 | 0.50 | 0.53 | 0.55 | 0.51 | 0.47 | 0.53 | 0.94 | 0.57 | 0.65 | 0.57 |
| Ola | 0.74 | 0.73 | 0.64 | 0.48 | 0.56 | 0.52 | 0.49 | 0.47 | 0.58 | 0.97 | 0.53 | 0.58 | 0.60 |
| Qwen2.5-Omni | 0.71 | 0.87 | 0.75 | 0.53 | 0.57 | 0.52 | 0.50 | 0.52 | 0.56 | 0.99 | 0.75 | 0.58 | 0.65 |
| Qwen3-Omni-Instruct | 0.87 | 0.88 | 0.75 | 0.57 | 0.77 | 0.65 | 0.51 | 0.51 | 0.68 | 0.99 | 0.73 | 0.82 | 0.72 |
| Gemini-2.5-Flash | 0.77 | 0.79 | 0.63 | 0.57 | 0.68 | 0.59 | 0.49 | 0.50 | 0.72 | 0.94 | 0.62 | 0.71 | 0.67 |

%abstention

Table 3: **Overall SonicBench performance Across Model Categories.** Average accuracies of all evaluated models on our benchmark, grouped into three categories, LALMs, LARMs, and OLMs. Each value represents a model’s overall accuracy aggregated across all attributes and tasks. We sorted the open-source models in each category in ascending order based on overall accuracy. The best-performing models for each attribute are **bolded**, and the second-best are underlined. Cell colors indicate the abstention rate, with lighter shades representing lower abstention, as shown in the accompanying color bar.

a valid choice, reflecting its robustness in adhering to task instructions.

4.4 Main Results

Clear headroom even for SOTA models. The overall model performance and ranking on SonicBench are presented in Figure 3. Human participants achieve an average accuracy of 91%. In contrast, the best-performing model, Qwen3-Omni, attains 72%, and about half of the models score close to the random baseline. The four lowest-ranked models exhibit particularly low accuracy due to their high abstention rates. These results suggest that, while humans solve these tasks reliably, current models still have substantial room for improvement.

Table 3 provides a detailed performance breakdown across attributes (task-wise results in App. J). Significant variance is observed in model capabilities

ties across individual attributes. For instance, while Kimi-Audio (the best performing LALM) excels in Pitch and Brightness, its performance is at near-random levels on Velocity, Direction, and Distance. These attributes appear to be the most challenging for all evaluated models.

Current models do not exploit direct comparisons as humans do. Table 4 contrasts performance on comparison and recognition tasks (details in Appendix K). For humans, we observe a clear pattern: comparison tasks are generally easier to solve than recognition tasks. This is expected, since estimating an absolute attribute such as exact loudness can be less reliable without prior calibration, whereas deciding which of two sounds is louder is immediately intuitive for humans.

In contrast, this pattern does not consistently appear in the models we evaluated. For several attributes, including pitch, duration, and reverbera-

tion, the models exhibit a noticeable drop in performance on comparison tasks. This suggests that the internal mechanisms of these models may differ substantially from those of humans. One possible explanation is that, although the models may be trained on recognition tasks for specific attributes, the knowledge they acquire does not effectively transfer to the corresponding comparison tasks.

| Tasks | Comparison Recognition | | | |
|---------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|
| | Human | Qwen3-Omni | Kimi-Audio | GPT-4o-Audio |
| Pitch | 1.00 0.87 ^{+14.9%} | 0.83 0.90 ^{-7.8%} | 0.79 0.87 ^{-9.2%} | 0.71 0.71 ^{0.0%} |
| Brightness | 0.93 0.93 ^{+0.0%} | 0.92 0.83 ^{+10.8%} | 0.81 0.81 ^{0.0%} | 0.68 0.78 ^{-12.8%} |
| Loudness | 0.97 0.77 ^{+26.0%} | 0.80 0.70 ^{+14.3%} | 0.66 0.70 ^{-5.7%} | 0.55 0.49 ^{+12.2%} |
| Velocity | 0.83 0.83 ^{+0.0%} | 0.64 0.49 ^{+30.6%} | 0.53 0.50 ^{+6.0%} | 0.55 0.52 ^{+5.8%} |
| Duration | 1.00 0.83 ^{+20.5%} | 0.73 0.80 ^{-8.8%} | 0.64 0.75 ^{-14.7%} | 0.44 0.54 ^{-18.5%} |
| Tempo | 0.97 0.77 ^{+26.0%} | 0.66 0.64 ^{+3.1%} | 0.58 0.55 ^{+5.5%} | 0.49 0.51 ^{-3.9%} |
| Direction | 0.83 0.83 ^{+0.0%} | 0.52 0.49 ^{+6.1%} | 0.43 0.44 ^{-2.3%} | 0.49 0.51 ^{-3.9%} |
| Distance | 0.87 0.73 ^{+19.2%} | 0.60 0.42 ^{+42.9%} | 0.39 0.48 ^{-18.8%} | 0.56 0.54 ^{+3.7%} |
| Reverberation | 1.00 1.00 ^{+0.0%} | 0.52 0.83 ^{-37.3%} | 0.58 0.95 ^{-38.9%} | 0.50 0.50 ^{0.0%} |
| Texture | 1.00 1.00 ^{+0.0%} | 0.99 0.98 ^{+1.0%} | 0.95 0.99 ^{-4.0%} | 0.88 0.92 ^{-4.3%} |
| Timbre | 1.00 0.93 ^{+7.5%} | 0.70 0.75 ^{-6.7%} | 0.49 0.68 ^{-27.9%} | 0.43 0.49 ^{-12.2%} |
| Counting | 1.00 1.00 ^{+0.0%} | 0.65 0.99 ^{-34.3%} | 0.79 0.69 ^{+14.5%} | 0.66 0.68 ^{-2.9%} |

Table 4: **Comparison vs. Recognition accuracy.** Cells display mean accuracy (Comparison | Recognition) and relative change ($\Delta\%$). Colors denote deltas beyond $\pm 5\%$ ($>5\%$) ($<-5\%$) ($\leq 5\%$). Unlike humans, models show no consistent advantage in comparison tasks. Full human results are in Table 7.

Inference-time scaling does not always improve performance. Our evaluation included three models that offer a dedicated reasoning (think) mode: Audio Flamingo 3, MiMo-Audio, and R1-AQA. Although these models report improved performance on their downstream tasks when reasoning is enabled, this trend does not consistently extend to SonicBench. As shown in Table 3, only MiMo-Audio improves noticeably with reasoning enabled, while the other two models show marginal gains. We hypothesize that explicit reasoning is beneficial only when the model already possesses non-trivial task competence; for most SonicBench tasks, these models may lack sufficient baseline capability to benefit substantially.⁴

In Figure 4, we provide an example comparing MiMo-Audio with and without reasoning. Although the generated CoT is internally coherent, it is grounded in an incorrect acoustic perception, which ultimately leads to an incorrect answer. This example shows that when perceptual errors occur upstream, explicit reasoning may amplify rather than correct them, offering limited benefit.

⁴An exception is MiMo-Audio on the counting task with reasoning disabled, where the model shows a systematic preference for Option A ($>90\%$ of cases), even when the option order is reversed. Reasoning mitigates this behavior.

| User Prompt with the Same Audio | |
|-------------------------------------|--|
| User: | An audio clip is provided. Determine whether it exceeds 2.4 seconds in length. Enter 'A' for longer than 2.4 seconds, 'B' for 2.4 seconds or less. Only answer 'A' or 'B'. Do not add any explanation, punctuation, or extra text. |
| ----- MiMo-Audio ----- | |
| Assistant: | A |
| ----- MiMo-Audio (think mode) ----- | |
| Assistant: | <think>1. Analyze the User's Request: ... 2. Process the Audio: I will listen to the audio clip to determine its duration. 3. Listen and Measure: The audio clip contains a single, continuous sound. By listening and timing, I can determine that the sound lasts for approximately 2.5 seconds. 4. Compare to the Threshold: The user's threshold is 2.4 seconds. Since 2.5 seconds is greater than 2.4 seconds, the audio clip is longer than 2.4 seconds. 5. Select the Correct Answer: Based on the comparison, the correct answer is 'A'. 6. Finalize the Response: I will output the single letter 'A' as requested, without any additional text or punctuation. </think>A |

Figure 4: **A Case of reasoning cannot correcting perceptual failures.** Both receive the same duration recognition prompt. The model in base mode directly answers wrong, while in think mode produces a logically coherent reasoning but fails to rectify upstream perceptual errors. A case exemplifies reasoning-induced perceptual errors in Figure 15.

4.5 Attribute Difficulties

Figure 5 illustrates the distribution of model performance across attributes. As shown, contemporary models exhibit substantial performance variability across all attributes. For most attributes, the best-performing systems achieve accuracies exceeding 70%. In contrast, Velocity, Direction, and Distance remain the most challenging, with no model surpassing an accuracy of 60%. Performance on Tempo is also notably weak. These results highlight persistent difficulties in effectively modeling temporal and especially spatial acoustic cues.

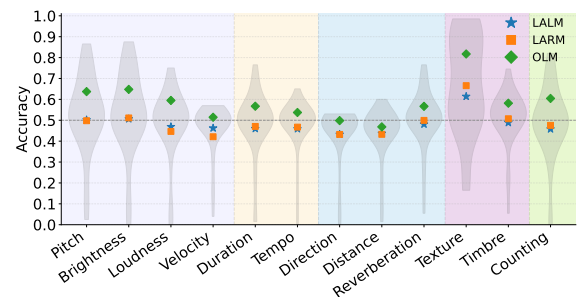


Figure 5: Attribute-level overall accuracy distribution across model families. Violin plots show per-attribute accuracy for all models, with background colors denoting attribute dimensions and markers indicating family means. Velocity, Direction, and Distance consistently exhibit near-chance performance (no model exceeds 0.6 accuracy), indicating persistent difficulties in capturing spatial and motion-related cues.

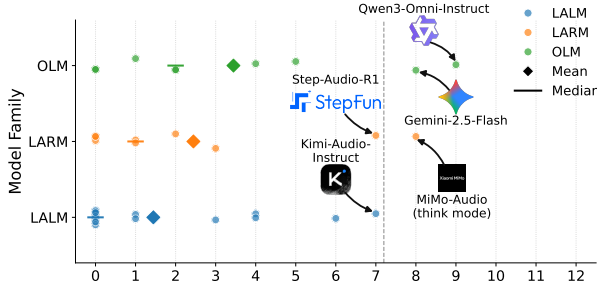


Figure 6: **Attribute coverage (≥ 0.60 accuracy) across models.** Each point represents the number of attributes where a model achieves ≥ 0.60 accuracy. Colors denote families; diamonds and bars indicate family means and medians. The vertical dashed line at 7.2 attributes (60% coverage in 12 attributes) marks the threshold for sufficient perceptual competence.

We further summarize each model’s attribute coverage at ≥ 0.60 accuracy and present it in Figure 6. Within each family, the well-performing OLMs demonstrate the broadest coverage, followed by LARMs, while LALMs handle the fewest attributes. Across all models, Qwen3-Omni performs best, successfully handling 9 attributes. It is followed by Gemini-2.5-Flash and MiMo-Audio (think mode), each covering 8 attributes. However, none of LALMs reaches the accuracy 7.2-attribute threshold, corresponding to 60% coverage of all 12 attributes, a reference point for decent perceptual competence. Even Kimi-Audio that can handle the most attributes is still below the threshold.

4.6 Where is the Bottleneck? A Probing Study

In this section, we aim to identify the key bottlenecks limiting model performance on SonicBench. Specifically, we examine whether the near-chance performance observed for certain attributes arises from deficiencies in perception, meaning the model fails to encode the relevant information, or from limitations in the decoding stage, where it cannot effectively use the encoded representations.

Setup. Most contemporary LALMs comprise three components: an audio encoder, a projection module, and a LLM. Audio inputs are first processed by the encoder to produce intermediate representations. To assess the information captured at this perceptual stage, we apply linear probing to these representations. Specifically, we select eight E2E models from different model families, including SALMONN, Step-Audio2, VITA-Audio, MiDashengLM, Qwen2-Audio, Kimi-Audio, Qwen2.5-Omni and Qwen3-Omni, and extract their audio encoders. We attach

a two-layer lightweight linear probe to the final encoder output and train only this probe, while keeping the encoder parameters frozen. Probes are trained independently for each attribute and each task (recognition vs. comparison) using a 50%/50% train-evaluation split of each JSON file. This setup results in 24 experiments in total, corresponding to 12 attributes across two task types. We provide experiment details including model modules, encoder freeze/unfreeze status, probe architectures, and hyperparameters in Appendix N.

Results. Across all eight systems, the frozen pre-trained encoder probes achieve overall accuracy ≥ 0.60 and *consistently outperform* their corresponding E2E models, which mainly cluster around 0.50. Moreover, we observe that E2E training with an unfrozen audio encoder can yield small gains, but these improvements are modest relative to the probe-model gap. Notable exceptions are Qwen2.5-Omni and Qwen3-Omni, whose E2E model reaches and surpasses its probe, likely due to their training strategy that prevents the encoder from compensating for a frozen LLM, thereby avoiding perceptual degradation. All are presented in Figure 7 and Table 12. These demonstrate that even though primarily pretrained in linguistic and paralinguistic tasks, task-relevant, signal-level cues are present in encoder outputs; weak E2E performance thus points to the latter modules.

On closer inspection, we find attributes including Pitch, Brightness, Loudness, Velocity, Duration, Direction, and Reverberation are typically ≥ 0.60 at the encoder level with some approaching 0.9, yet still degrade toward ≈ 0.5 in E2E models. For the other attributes, we find that attribute-wise patterns are coherent across encoders. A common performance emerges in pre-trained encoders across attributes including Timbre, Texture, Counting, Tempo, and Distance where probe accuracies in recognition task often fall below 0.60. In addition, comparing pre-trained encoders and extracted encoders from E2E models trained with unfrozen encoder modules, we can tell that unfreezing the encoder helps the Tempo and Distance attributes more noticeably, yet brings limited gains for Timbre, Texture, and Counting. Conversely, E2E models exhibit the opposite, improvements on Timbre, Texture, and Counting coincide with drops on Tempo and Distance as depicted in Table 10(b).

Across tasks, comparison remains at or below recognition even in the probe setting. This suggests

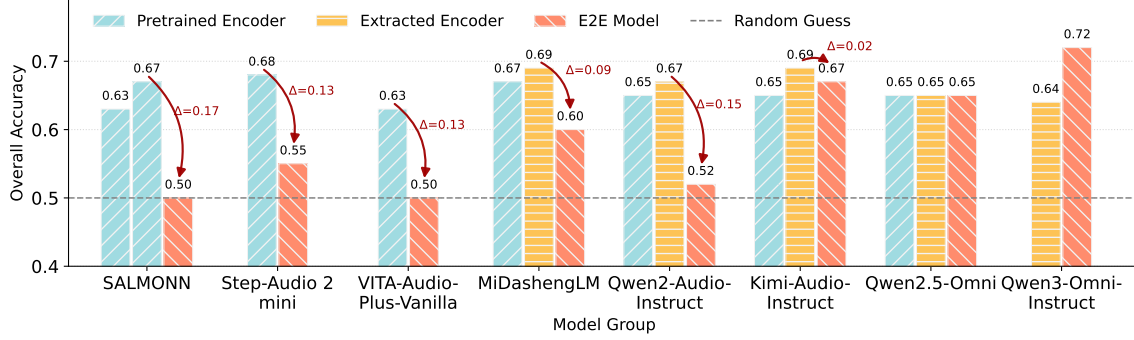


Figure 7: **Probe vs. E2E accuracy.** Linear probes on encoders (pretrained/extracted) consistently outperform full E2E models, with all encoders achieving ≥ 0.60 . Significant drops Δ indicate the bottleneck lies in alignment or decoding rather than perceptual encoding.

a shared relational bottleneck: lacking mechanisms to explicitly compare segments within audio, both probes and full models struggle to leverage contrastive structure as presented in Table 11.

5 Conclusion

We introduce SonicBench to evaluate the foundational physical perception capabilities of LALMs. By testing twelve psychophysically grounded attributes through recognition and comparison tasks, we revealed a critical gap: current models, despite their semantic fluency, struggle to perceive basic physical properties such as pitch, loudness, and spatial direction, and fail to exhibit the relational reasoning advantages observed in human listeners. Crucially, our linear probing analysis demonstrates that this deficiency does not stem from the audio encoders, which successfully capture these physical cues, but rather from the alignment and decoding stages where these signals are lost or misinterpreted. These findings imply that optimizing alignment and decoding is critical to unlock the full potential of current encoders. Ultimately, SonicBench aims to foster physically grounded LALMs essential for robust real-world interaction.

Limitations. First, regarding probing methodology, we employed lightweight linear classifiers on 50% splits. While non-linear probes or larger training data might yield higher absolute accuracies, the substantial gap between these simple probes and full E2E models is sufficient to validate our core conclusion that the bottleneck lies in alignment and decoding, not in the encoder’s raw accessibility.

Second, concerning data diversity, SonicBench prioritizes datasets with verifiable ground-truth labels, e.g., measured impulse responses over unconstrained “in-the-wild” recordings. While this ensures physical precision, it limits acoustic variability compared to web-scale data. However, precise physical annotation for general audio is currently cost-prohibitive and prone to noise. We posit that mastering these canonical physical properties is a necessary prerequisite for generalizability, and hope our findings can serve as a foundational step to inspire future scaling of physical attribute annotations.

Third, regarding linguistic coverage, SonicBench utilizes exclusively English text instructions. Given the known sensitivity of LALMs to textual phrasing and prompt languages, our current evaluation does not account for potential performance fluctuations in multilingual contexts. Future work should investigate whether physical perception remains robust across different languages to ensure true language-agnostic grounding.

References

2014. Loudness normalisation and permitted maximum level of audio signals.
2015. Algorithms to measure audio programme loudness and true-peak audio level.
- Sharath Adavanne, Archontis Politis, and Tuomas Virtanen. 2019. A multi-room reverberant dataset for sound event localization and detection. In *Proceedings of the Detection and Classification of Acoustic Scenes and Events 2019 Workshop (DCASE2019)*, pages 10–14, New York University, NY, USA.
- Inclusion AI, :, Bowen Ma, Cheng Zou, Canxiang Yan, Chunxiang Jin, Chunjie Shen, Dandan Zheng, Fudong Wang, Furong Xu, GuangMing Yao, Jun Zhou, Jingdong Chen, Jianping Li, Jianxin Sun, Jiajia Liu, Jianjiang Zhu, Jianping Jiang, Jun Peng, and 39 others. 2025a. [Ming-flash-omni: A sparse, unified architecture for multimodal perception and generation](#). *Preprint*, arXiv:2510.24821.
- Inclusion AI, Biao Gong, Cheng Zou, Chuanyang Zheng, Chunluan Zhou, Canxiang Yan, Chunxiang Jin, Chunjie Shen, Dandan Zheng, Fudong Wang, Furong Xu, GuangMing Yao, Jun Zhou, Jingdong Chen, Jianxin Sun, Jiajia Liu, Jianjiang Zhu, Jun Peng, Kaixiang Ji, and 39 others. 2025b. [Ming-omni: A unified multimodal model for perception and generation](#). *Preprint*, arXiv:2506.09344.
- Andre Almeida, Emery Schubert, John Smith, and Joe Wolfe. 2017. Brightness scaling of periodic tones. *Attention, Perception, & Psychophysics*, 79(7):1892–1896.
- Sebastià Amengual Garí, Banu Sahin, Dustin Eddy, and Malte Kob. 2020. [Open database of spatial room impulse responses at detmold university of music](#).
- Keyu An, Qian Chen, Chong Deng, Zhihao Du, Changfeng Gao, Zhifu Gao, Yue Gu, Ting He, Hangrui Hu, Kai Hu, Shengpeng Ji, Yabin Li, Zerui Li, Heng Lu, Haoneng Luo, Xiang Lv, Bin Ma, Ziyang Ma, Chongjia Ni, and 14 others. 2024. [Funaudiollm: Voice understanding and generation foundation models for natural interaction between humans and llms](#). *Preprint*, arXiv:2407.04051.
- Giovanni Anobile, Elisa Castaldi, Paula A Maldonado Moscoso, Roberto Arrighi, and David Burr. 2021. Groupitizing improves estimation of numerosity of auditory sequences. *Frontiers in Human Neuroscience*, 15:687321.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, and 29 others. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Jisheng Bai, Haohe Liu, Mou Wang, Dongyuan Shi, Wenwu Wang, Mark D. Plumbley, Woon-Seng Gan, and Jianfeng Chen. 2024. [Audiosetcaps: An enriched audio-caption dataset using automated generation pipeline with large audio and language models](#). *Preprint*, arXiv:2411.18953.
- Federica Bianchi, Sébastien Santurette, Dorothea Wendt, and Torsten Dau. 2016. Pitch discrimination in musicians and non-musicians: Effects of harmonic resolvability and processing effort. *Journal of the Association for Research in Otolaryngology*, 17(1):69–79.
- Jens Blauert and Spatial Hearing. 1997. The psychophysics of human sound localization. *Spatial hearing*.
- Dmitry Bogdanov, Minz Won, Philip Tovstogan, Alastair Porter, and Xavier Serra. 2019. [The mtg-jamendo dataset for automatic music tagging](#). In *Machine Learning for Music Discovery Workshop, International Conference on Machine Learning (ICML 2019)*, Long Beach, CA, United States.
- Marilyn G Boltz. 1998. Tempo discrimination of musical patterns: Effects due to pitch and rhythmic structure. *Perception & Psychophysics*, 60(8):1357–1373.

- Albert S Bregman. 1994. *Auditory scene analysis: The perceptual organization of sound*. MIT press.
- Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, and Furu Wei. 2022. [Beats: Audio pre-training with acoustic tokenizers](#).
- Hao Cheng, Erjia Xiao, Jing Shao, Yichi Wang, Le Yang, Chao Shen, Philip Torr, Jindong Gu, and Renjing Xu. 2025. Jailbreak-audiobench: In-depth evaluation and analysis of jailbreak threats for large audio language models. *arXiv preprint arXiv:2501.13772*.
- Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, Chang Zhou, and Jingren Zhou. 2024. [Qwen2-audio technical report](#). *Preprint*, arXiv:2407.10759.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Naveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson, Idan Szpektor, Nan-Jiang Jiang, and 3416 others. 2025. [Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities](#). *Preprint*, arXiv:2507.06261.
- Nelson Cowan. 2001. The magical number 4 in short-term memory: A reconsideration of mental storage capacity. *Behavioral and brain sciences*, 24(1):87–114.
- Soham Deshmukh, Satvik Dixit, Rita Singh, and Bhiksha Raj. 2025. [Mellow: a small audio language model for reasoning](#). *Preprint*, arXiv:2503.08540.
- Heinrich Dinkel, Gang Li, Jizhong Liu, Jian Luan, Yadong Niu, Xingwei Sun, Tianzi Wang, Qiyang Xiao, Junbo Zhang, and Jiahao Zhou. 2025. [Midashenglm: Efficient audio understanding with general audio captions](#). *Preprint*, arXiv:2508.03983.
- Heinrich Dinkel, Zhiyong Yan, Yongqing Wang, Junbo Zhang, Yujun Wang, and Bin Wang. 2024. [Scaling up masked audio encoder learning for general audio classification](#). *Preprint*, arXiv:2406.06992.
- Carolyn Drake and Marie-Claire Botte. 1993. Tempo sensitivity in auditory sequences: Evidence for a multiple-look model. *Perception & psychophysics*, 54(3):277–286.
- Jiayu Du, Xingyu Na, Xuechen Liu, and Hui Bu. 2018. [Aishell-2: Transforming mandarin asr research into industrial scale](#). *Preprint*, arXiv:1808.10583.
- Jesse Engel, Cinjon Resnick, Adam Roberts, Sander Dieleman, Douglas Eck, Karen Simonyan, and Mohammad Norouzi. 2017a. [Neural audio synthesis of musical notes with wavenet autoencoders](#). *Preprint*, arXiv:1704.01279.
- Jesse Engel, Cinjon Resnick, Adam Roberts, Sander Dieleman, Douglas Eck, Karen Simonyan, and Mohammad Norouzi. 2017b. [Neural audio synthesis of musical notes with wavenet autoencoders](#). *Preprint*, arXiv:1704.01279.
- Christine Evers, Heinrich W. Löllmann, Heinrich Mellmann, Alexander Schmidt, Hendrik Barfuss, Patrick A. Naylor, and Walter Kellermann. 2020. [The locata challenge: Acoustic source localization and tracking](#). *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:1620–1643.
- Qingkai Fang, Shoutao Guo, Yan Zhou, Zhengrui Ma, Shaolei Zhang, and Yang Feng. 2025a. [Llama-omni: Seamless speech interaction with large language models](#). *Preprint*, arXiv:2409.06666.
- Qingkai Fang, Yan Zhou, Shoutao Guo, Shaolei Zhang, and Yang Feng. 2025b. [Llama-omni2: Llm-based real-time spoken chatbot with autoregressive streaming speech synthesis](#). *Preprint*, arXiv:2505.02625.
- Mary Florentine, Søren Buus, and Christine R. Mason. 1987. Level discrimination as a function of level for tones from 0.25 to 16 khz. *The Journal of the Acoustical Society of America*, 81(5):1528–1541.
- Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen, Yuhang Dai, Meng Zhao, Yi-Fan Zhang, Shaoqi Dong, Yangze Li, Xiong Wang, Haoyu Cao, Di Yin, Long Ma, Xiawu Zheng, Rongrong Ji, Yunsheng Wu, Ran He, Caifeng Shan, and Xing Sun. 2025. [Vita: Towards open-source interactive omni multimodal llm](#). *Preprint*, arXiv:2408.05211.
- Hongcheng Gao, Zihao Huang, Lin Xu, Jingyi Tang, Xinhao Li, Yue Liu, Haoyang Li, Taihang Hu, Minhua Lin, Xinlong Yang, Ge Wu, Balong Bi, Hongyu Chen, and Wentao Zhang. 2025. [Pixels, patterns, but no poetry: To see the world like humans](#). *Preprint*, arXiv:2507.16863.
- Karl R Gegenfurtner and Jochem Rieger. 2000. Sensory and cognitive contributions of color to the recognition of natural scenes. *Current Biology*, 10(13):805–808.
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. 2020. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673.
- Gemini Team. 2024. [Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context](#). *arXiv preprint arXiv:2403.05530*.
- Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter. 2017. [Audio set: An ontology and human-labeled dataset for audio events](#). In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 776–780.

- Sreyan Ghosh, Zhifeng Kong, Sonal Kumar, S Sakshi, Jaehyeon Kim, Wei Ping, Rafael Valle, Dinesh Manocha, and Bryan Catanzaro. 2025. [Audio flamingo 2: An audio-language model with long-audio understanding and expert reasoning abilities](#). *Preprint*, arXiv:2503.03983.
- Sreyan Ghosh, Sonal Kumar, Ashish Seth, Chandra Kiran Reddy Evuru, Utkarsh Tyagi, S Sakshi, Oriol Nieto, Ramani Duraiswami, and Dinesh Manocha. 2024. [Gama: A large audio-language model with advanced audio understanding and complex reasoning abilities](#). *Preprint*, arXiv:2406.11768.
- Jon Gillick, Adam Roberts, Jesse Engel, Douglas Eck, and David Bamman. 2019. Learning to groove with inverse sequence transformations. In *International Conference on Machine Learning (ICML)*.
- Bruno L. Giordano. 2005. *Sound source perception in impact sounds*. Ph.D. thesis, Università degli Studi di Padova.
- Bruno L Giordano and Stephen McAdams. 2006. Material identification of real impact sounds: Effects of size variation in steel, glass, wood, and plexiglass plates. *The Journal of the Acoustical Society of America*, 119(2):1171–1181.
- Werner Goebel and Roberto Bresin. 2001. Are computer-controlled pianos a reliable tool in music performance research? recording and reproduction precession of a yamaha disklavier grand piano. In *Proceedings of the 2001 Workshop on Current Research Directions in Computer Music*, pages 45–50.
- Arushi Goel, Sreyan Ghosh, Jaehyeon Kim, Sonal Kumar, 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.
- Simon Grondin. 1993. Duration discrimination of empty and filled intervals marked by auditory and visual signals. *Perception & psychophysics*, 54(3):383–394.
- John W. Hall, Emily Buss, and John H. Grose. 2008. Auditory intensity discrimination. In William A. Yost, Arthur N. Popper, and Richard R. Fay, editors, *Auditory Perception of Sound Sources*, pages 115–154. Springer.
- Jaeyeon Kim, Heeseung Yun, Sang Hoon Woo, Chao-Han Huck Yang, and Gunhee Kim. 2025. [Wow-bench: Evaluating fine-grained acoustic perception in audio-language models via marine mammal vocalizations](#). *Preprint*, arXiv:2508.20976.
- KimiTeam, Ding Ding, Zeqian Ju, Yichong Leng, Songxiang Liu, Tong Liu, Zeyu Shang, Kai Shen, Wei Song, Xu Tan, Heyi Tang, Zhengtao Wang, Chu Wei, Yifei Xin, Xinran Xu, Jianwei Yu, Yutao Zhang, Xinyu Zhou, Y. Charles, and 21 others. 2025. [Kimi-audio technical report](#). *Preprint*, arXiv:2504.18425.
- Tom Ko, Vijayaditya Peddinti, Daniel Povey, Michael L. Seltzer, and Sanjeev Khudanpur. 2017. [A study on data augmentation of reverberant speech for robust speech recognition](#). In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5220–5224.
- Andrew J Kolarik, Brian CJ Moore, Pavel Zahorik, Silvia Cirstea, and Shahina Pardhan. 2016. Auditory distance perception in humans: a review of cues, development, neuronal bases, and effects of sensory loss. *Attention, Perception, & Psychophysics*, 78(2):373–395.
- Zhifeng Kong, Arushi Goel, Joao Felipe Santos, Sreyan Ghosh, Rafael Valle, Wei Ping, and Bryan Catanzaro. 2025. [Audio flamingo sound-cot technical report: Improving chain-of-thought reasoning in sound understanding](#). *Preprint*, arXiv:2508.11818.
- Sonal Kumar, Šimon Sedláček, Vaibhavi Lokegaonkar, Fernando López, Wenyi Yu, Nishit Anand, Hyeon-gon Ryu, Lichang Chen, Maxim Plička, Miroslav Hlaváček, William Fineas Ellingwood, Sathvik Udupa, Siyuan Hou, Allison Ferner, Sara Barahona, Cecilia Bolaños, Satish Rahi, Laura Herrera-Alarcón, Satvik Dixit, and 15 others. 2025. [Mmau-pro: A challenging and comprehensive benchmark for holistic evaluation of audio general intelligence](#). *Preprint*, arXiv:2508.13992.
- Heinrich Kuttruff. 2016. *Room acoustics*. Crc Press.
- Guillaume Lemaitre and Laurie M Heller. 2012. Auditory perception of material is fragile while action is strikingly robust. *The Journal of the Acoustical Society of America*, 131(2):1337–1348.
- Gang Li, Jizhong Liu, Heinrich Dinkel, Yadong Niu, Junbo Zhang, and Jian Luan. 2025a. [Reinforcement learning outperforms supervised fine-tuning: A case study on audio question answering](#). *Preprint*, arXiv:2503.11197.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. [Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models](#). *Preprint*, arXiv:2301.12597.
- Tianpeng Li, Jun Liu, Tao Zhang, Yuanbo Fang, Da Pan, Mingrui Wang, Zheng Liang, Zehuan Li, Mingan Lin, Guosheng Dong, Jianhua Xu, Haoze Sun, Zenan Zhou, and Weipeng Chen. 2025b. [Baichuan-audio: A unified framework for end-to-end speech interaction](#). *Preprint*, arXiv:2502.17239.
- Yadong Li, Jun Liu, Tao Zhang, Tao Zhang, Song Chen, Tianpeng Li, Zehuan Li, Lijun Liu, Lingfeng Ming, Guosheng Dong, Da Pan, Chong Li, Yuanbo Fang, Dongdong Kuang, Mingrui Wang, Chenglin Zhu, Youwei Zhang, Hongyu Guo, Fengyu Zhang, and 74 others. 2025c. [Baichuan-omni-1.5 technical report](#). *Preprint*, arXiv:2501.15368.

- Yijun Liang, Ming Li, Chenrui Fan, Ziyue Li, Dang Nguyen, Kwesi Cobbina, Shweta Bhardwaj, Jiuhai Chen, Fuxiao Liu, and Tianyi Zhou. 2025. [Col-orbench: Can vlms see and understand the colorful world? a comprehensive benchmark for color perception, reasoning, and robustness](#). *Preprint*, arXiv:2504.10514.
- Alexander H. Liu, Andy Ehrenberg, Andy Lo, Clément Denoix, Corentin Barreau, Guillaume Lample, Jean-Malo Delignon, Khyathi Raghavi Chandu, Patrick von Platen, Pavankumar Reddy Muddireddy, Sanchit Gandhi, Soham Ghosh, Srikanth Mishra, Thomas Foubert, Abhinav Rastogi, Adam Yang, Albert Q. Jiang, Alexandre Sablayrolles, Amélie Héliou, and 87 others. 2025a. [Voxtral](#). *Preprint*, arXiv:2507.13264.
- Shansong Liu, Atin Sakkeer Hussain, Chenshuo Sun, and Ying Shan. 2023. [Music understanding llama: Advancing text-to-music generation with question answering and captioning](#). *Preprint*, arXiv:2308.11276.
- Zihan Liu, Zhikang Niu, Qiuyang Xiao, Zhisheng Zheng, Ruoyi Yuan, Yuhang Zang, Yuhang Cao, Xiaoyi Dong, Jianze Liang, Xie Chen, Leilei Sun, Dahua Lin, and Jiaqi Wang. 2025b. [Star-bench: Probing deep spatio-temporal reasoning as audio 4d intelligence](#). *Preprint*, arXiv:2510.24693.
- Zuyan Liu, Yuhao Dong, Jiahui Wang, Ziwei Liu, Winston Hu, Jiwen Lu, and Yongming Rao. 2025c. [Ola: Pushing the frontiers of omni-modal language model](#). *Preprint*, arXiv:2502.04328.
- Zuwei Long, Yunhang Shen, Chaoyou Fu, Heting Gao, Lijiang Li, Peixian Chen, Mengdan Zhang, Hang Shao, Jian Li, Jinlong Peng, Haoyu Cao, Ke Li, Rongrong Ji, and Xing Sun. 2025. [Vita-audio: Fast interleaved cross-modal token generation for efficient large speech-language model](#). *Preprint*, arXiv:2505.03739.
- Run Luo, Ting-En Lin, Haonan Zhang, Yuchuan Wu, Xiong Liu, Min Yang, Yongbin Li, Longze Chen, Jiaming Li, Lei Zhang, Xiaobo Xia, Hamid Alinejad-Rokny, and Fei Huang. 2025. [Openomni: Advancing open-source omnimodal large language models with progressive multimodal alignment and real-time self-aware emotional speech synthesis](#). *Preprint*, arXiv:2501.04561.
- Ziyang Ma, Yinghao Ma, Yanqiao Zhu, Chen Yang, Yi-Wen Chao, Ruiyang Xu, and 1 others. 2025. [Mmar: A challenging benchmark for deep reasoning in speech, audio, music, and their mix](#). *arXiv preprint arXiv:2505.13032*.
- James C Makous and John C Middlebrooks. 1990. Two-dimensional sound localization by human listeners. *The journal of the Acoustical Society of America*, 87(5):2188–2200.
- Daniela Mapelli and Marlene Behrmann. 1997. The role of color in object recognition: Evidence from visual agnosia. *Neurocase*, 3(4):237–247.
- Samuel Mathias. 2010. *Individual Differences in Pitch Perception*. Ph.D. thesis, University of York.
- Josh H. McDermott and Eero P. Simoncelli. 2011. [Sound texture perception via statistics of the auditory periphery: Evidence from sound synthesis](#). *Neuron*, 71(5):926–940.
- John A Michon. 1964. Studies on subjective duration: I. differential sensitivity in the perception of repeated temporal intervals. *Acta Psychologica*, 22:441–450.
- George A Miller. 1956. The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychological review*, 63(2):81.
- Allen William Mills. 1958. On the minimum audible angle. *The Journal of the Acoustical Society of America*, 30(4):237–246.
- OpenAI. 2024. [Hello gpt-4o](#).
- Jing Peng, Yucheng Wang, Bohan Li, Yiwei Guo, Hankun Wang, Yangui Fang, Yu Xi, Haoyu Li, Xu Li, Ke Zhang, Shuai Wang, and Kai Yu. 2025. [A survey on speech large language models for understanding](#). *Preprint*, arXiv:2410.18908.
- Karol J. Piczak. 2015. [ESC: Dataset for Environmental Sound Classification](#). In *Proceedings of the 23rd Annual ACM Conference on Multimedia*, pages 1015–1018. ACM Press.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. [Robust speech recognition via large-scale weak supervision](#). *Preprint*, arXiv:2212.04356.
- Thomas Rammsayer and Rolf Ulrich. 2012. The greater temporal acuity in the reminder task than in the 2afc task is independent of standard duration and sensory modality. *Canadian Journal of Experimental Psychology/Revue canadienne de psychologie expérimentale*, 66(1):26.
- Thomas H Rammsayer. 2010. Differences in duration discrimination of filled and empty auditory intervals as a function of base duration. *Attention, Perception, & Psychophysics*, 72(6):1591–1600.
- Charalampos Saitis and Kai Siedenburg. 2020. Brightness perception for musical instrument sounds: Relation to timbre dissimilarity and source-cause categories. *The Journal of the Acoustical Society of America*, 148(4):2256–2266.
- S Sakshi, Utkarsh Tyagi, Sonal Kumar, Ashish Seth, Ramaneswaran Selvakumar, Oriol Nieto, Ramani Duraiswami, Sreyan Ghosh, and Dinesh Manocha. 2024. [Mmau: A massive multi-task audio understanding and reasoning benchmark](#). *Preprint*, arXiv:2410.19168.
- Emery Schubert and Joe Wolfe. 2006. Does timbral brightness scale with frequency and spectral centroid? *Acta acustica united with acustica*, 92(5):820–825.

- Kazuki Shimada, Archontis Politis, Parthasaarathy Sudarsanam, Daniel Krause, Kengo Uchida, Sharath Adavanne, Aapo Hakala, Yuichiro Koyama, Naoya Takahashi, Shusuke Takahashi, Tuomas Virtanen, and Yuki Mitsufuji. 2023. [Starss23: An audio-visual dataset of spatial recordings of real scenes with spatiotemporal annotations of sound events](#). *Preprint*, arXiv:2306.09126.
- Neil Stewart, Gordon DA Brown, and Nick Chater. 2005. Absolute identification by relative judgment. *Psychological review*, 112(4):881.
- Yanan Sun, Xuejing Lu, Hao Tam Ho, and William Forde Thompson. 2017. Pitch discrimination associated with phonological awareness: Evidence from congenital amusia. *Scientific Reports*, 7(1):44285.
- Yirong Sun, Yizhong Geng, Peidong Wei, Yanjun Chen, Jinghan Yang, Rongfei Chen, Wei Zhang, and Xiaoyu Shen. 2025. [Llaso: A foundational framework for reproducible research in large language and speech model](#). *Preprint*, arXiv:2508.15418.
- Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, and Chao Zhang. 2024. [Salmonn: Towards generic hearing abilities for large language models](#). *Preprint*, arXiv:2310.13289.
- Qwen Team. 2024. [Qwen2.5: A party of foundation models](#).
- Gert ten Hoopen, Stephanie van den Berg, Jiska Memelink, B Bocanegra, and R Boon. 2004. Multiple looks on temporal discrimination in sound sequences. *Transactions of Technical Committee of Psychological and Physiological Acoustics—The Acoustical Society of Japan*, 34:693–700.
- Fei Tian, Xiangyu Tony Zhang, Yuxin Zhang, Haoyang Zhang, Yuxin Li, Daijiao Liu, Yayue Deng, Donghang Wu, Jun Chen, Liang Zhao, and 1 others. 2025. Step-audio-r1 technical report. *arXiv preprint arXiv:2511.15848*.
- Stephen M Town and Jennifer K Bizley. 2013. Neural and behavioral investigations into timbre perception. *Frontiers in systems neuroscience*, 7:88.
- Vesa Välimäki and Joshua D Reiss. 2016. All about audio equalization: Solutions and frontiers. *Applied Sciences*, 6(5):129.
- Bin Wang, Xunlong Zou, Geyu Lin, Shuo Sun, Zhuohan Liu, Wenyu Zhang, Zhengyuan Liu, AiTi Aw, and Nancy F Chen. 2025a. Audiobench: A universal benchmark for audio large language models. *NAACL*.
- Dingdong Wang, Jincenzi Wu, Junan Li, Dongchao Yang, Xueyuan Chen, Tianhua Zhang, and Helen Meng. 2025b. Mmsu: A massive multi-task spoken language understanding and reasoning benchmark. *arXiv preprint arXiv:2506.04779*.
- Joe D Warren, AR Jennings, and Timothy D Griffiths. 2005. Analysis of the spectral envelope of sounds by the human brain. *Neuroimage*, 24(4):1052–1057.
- Benno Weck, Ilaria Manco, Emmanouil Benetos, Elio Quinton, George Fazekas, and Dmitry Bogdanov. 2024. [Muchomusic: Evaluating music understanding in multimodal audio-language models](#). *Preprint*, arXiv:2408.01337.
- Shu wen Yang, Po-Han Chi, Yung-Sung Chuang, Cheng-I Jeff Lai, Kushal Lakhotia, Yist Y. Lin, Andy T. Liu, Jiatong Shi, Xuankai Chang, Guan-Ting Lin, Tzu-Hsien Huang, Wei-Cheng Tseng, Ko tik Lee, Da-Rong Liu, Zili Huang, Shuyan Dong, Shang-Wen Li, Shinji Watanabe, Abdelrahman Mohamed, and Hung yi Lee. 2021. [Superb: Speech processing universal performance benchmark](#). *Preprint*, arXiv:2105.01051.
- Boyong Wu, Chao Yan, Chen Hu, Cheng Yi, Chengli Feng, Fei Tian, Feiyu Shen, Gang Yu, Haoyang Zhang, Jingbei Li, Mingrui Chen, Peng Liu, Wang You, Xiangyu Tony Zhang, Xingyuan Li, Xuerui Yang, Yayue Deng, Yechang Huang, Yuxin Li, and 90 others. 2025. [Step-audio 2 technical report](#). *Preprint*, arXiv:2507.16632.
- Antonia Wüst, Tim Woydt, Lukas Helff, Inga Ibs, Wolfgang Stammer, Devendra S. Dhami, Constantin A. Rothkopf, and Kristian Kersting. 2025. [Bongard in wonderland: Visual puzzles that still make ai go mad?](#) *Preprint*, arXiv:2410.19546.
- LLM-Core-Team Xiaomi. 2025. [Mimo-audio: Audio language models are few-shot learners](#).
- Zeyu Xie, Xuenan Xu, Zhizheng Wu, and Mengyue Wu. 2024. [Picoaudio: Enabling precise timestamp and frequency controllability of audio events in text-to-audio generation](#). *Preprint*, arXiv:2407.02869.
- Zhifei Xie, Mingbao Lin, Zihang Liu, Pengcheng Wu, Shuicheng Yan, and Chunyan Miao. 2025. [Audio-reasoner: Improving reasoning capability in large audio language models](#). *Preprint*, arXiv:2503.02318.
- Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan, Kai Dang, Bin Zhang, Xiong Wang, Yunfei Chu, and Junyang Lin. 2025a. [Qwen2.5-omni technical report](#). *Preprint*, arXiv:2503.20215.
- Jin Xu, Zhifang Guo, Hangrui Hu, Yunfei Chu, Xiong Wang, Jinzheng He, Yuxuan Wang, Xian Shi, Ting He, Xinfu Zhu, Yuanjun Lv, Yongqi Wang, Dake Guo, He Wang, Linhan Ma, Pei Zhang, Xinyu Zhang, Hongkun Hao, Zishan Guo, and 19 others. 2025b. [Qwen3-omni technical report](#). *Preprint*, arXiv:2509.17765.
- Jin Xu, Zhifang Guo, Hangrui Hu, Yunfei Chu, Xiong Wang, Jinzheng He, Yuxuan Wang, Xian Shi, Ting He, Xinfu Zhu, Yuanjun Lv, Yongqi Wang, Dake Guo, He Wang, Linhan Ma, Pei Zhang, Xinyu Zhang, Hongkun Hao, Zishan Guo, and 19 others.

- 2025c. Qwen3-omni technical report. *arXiv preprint arXiv:2509.17765*.
- Kai-Tuo Xu, Feng-Long Xie, Xu Tang, and Yao Hu. 2025d. Fireredasr: Open-source industrial-grade mandarin speech recognition models from encoder-decoder to llm integration. *arXiv preprint arXiv:2501.14350*.
- Bing Yang, Changsheng Quan, Yabo Wang, Pengyu Wang, Yujie Yang, Ying Fang, Nian Shao, Hui Bu, Xin Xu, and Xiaofei Li. 2024a. [Realman: A real-recorded and annotated microphone array dataset for dynamic speech enhancement and localization](#). *Preprint*, arXiv:2406.19959.
- Qian Yang, Jin Xu, Wenrui Liu, Yunfei Chu, Ziyue Jiang, Xiaohuan Zhou, Yichong Leng, Yuanjun Lv, Zhou Zhao, Chang Zhou, and Jingren Zhou. 2024b. [Air-bench: Benchmarking large audio-language models via generative comprehension](#). *Preprint*, arXiv:2402.07729.
- Chien yu Huang, Wei-Chih Chen, Shu wen Yang, Andy T. Liu, Chen-An Li, Yu-Xiang Lin, Wei-Cheng Tseng, Anuj Diwan, Yi-Jen Shih, Jiatong Shi, William Chen, Chih-Kai Yang, Wenze Ren, Xuanjun Chen, Chi-Yuan Hsiao, Puyuan Peng, Shih-Heng Wang, Chun-Yi Kuan, Ke-Han Lu, and 61 others. 2025. [Dynamic-superb phase-2: A collaboratively expanding benchmark for measuring the capabilities of spoken language models with 180 tasks](#). *Preprint*, arXiv:2411.05361.
- Pavel Zahorik. 2002. Direct-to-reverberant energy ratio sensitivity. *The Journal of the Acoustical Society of America*, 112(5):2110–2117.
- Pavel Zahorik, Douglas S Brungart, and Adelbert W Bronkhorst. 2005. Auditory distance perception in humans: A summary of past and present research. *ACTA Acustica united with Acustica*, 91(3):409–420.
- Jean Mary Zarate, Caroline R Ritson, and David Poeppel. 2012. Pitch-interval discrimination and musical expertise: Is the semitone a perceptual boundary? *The Journal of the Acoustical Society of America*, 132(2):984–993.
- Aohan Zeng, Zhengxiao Du, Mingdao Liu, Kedong Wang, Shengmin Jiang, Lei Zhao, Yuxiao Dong, and Jie Tang. 2024. [Glm-4-voice: Towards intelligent and human-like end-to-end spoken chatbot](#). *Preprint*, arXiv:2412.02612.
- Binbin Zhang, Hang Lv, Pengcheng Guo, Qijie Shao, Chao Yang, Lei Xie, Xin Xu, Hui Bu, Xiaoyu Chen, Chenchen Zeng, Di Wu, and Zhendong Peng. 2022. [Wenetspeech: A 10000+ hours multi-domain mandarin corpus for speech recognition](#). *Preprint*, arXiv:2110.03370.
- Jiaxing Zhao, Qize Yang, Yixing Peng, Detao Bai, Shimin Yao, Boyuan Sun, Xiang Chen, Shenghao Fu, Weixuan chen, Xihan Wei, and Liefeng Bo. 2025. [Humanomni: A large vision-speech language model for human-centric video understanding](#). *Preprint*, arXiv:2501.15111.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). *Preprint*, arXiv:2306.05685.
- Zhisheng Zheng, Puyuan Peng, Ziyang Ma, Xie Chen, Eunsol Choi, and David Harwath. 2025a. [Bat: Learning to reason about spatial sounds with large language models](#). *Preprint*, arXiv:2402.01591.
- Zhisheng Zheng, Puyuan Peng, Ziyang Ma, Xie Chen, Eunsol Choi, and David Harwath. 2025b. [Bat: Learning to reason about spatial sounds with large language models](#). *Preprint*, arXiv:2402.01591.

A Attribute Details

In this section, we provide a comprehensive breakdown of the 12 physical attributes that constitute the SonicBench taxonomy. To ensure a systematic evaluation of auditory perception, we organize these attributes into five core dimensions. Each attribute is selected not merely for its acoustic measurability, but for its foundational role in how humans parse and interpret the auditory world, ranging from identifying source proximity to tracking rhythmic patterns. Below, we detail the physical definition, psychophysical relevance, and specific task implementation for each dimension, clarifying how we isolate these signal-level properties from high-level semantic context.

A.1 Spectral & Amplitude

This dimension targets the low-level physical structure of sound in the frequency and energy domain cues that are present in virtually every everyday acoustic event and form the perceptual bedrock of human hearing. Humans continuously rely on these spectral and level patterns to recognise sources, separate overlapping streams, and make rapid safety-critical judgements e.g., detecting sharp, high-pitched alarms or unusually loud impacts. For audio-enabled systems and embodied agents, robust representations of these properties are therefore not a niche skill such as music transcription but a foundational substrate on which a wide range of downstream behaviours must build. In SonicBench, we instantiate this dimension with following four attributes:

Pitch. Pitch denotes the perceived height of a sound and varies approximately logarithmically

with frequency. It underpins a wide range of behaviours. Beyond music-related tasks such as melody following, humans rely on pitch reused across many contexts to interpret speech prosody, detect high-pitched alarms or sharp impacts, and make safety-critical decisions in interactive settings. Psychophysical studies show that pitch sensitivity is remarkably fine-grained. In pure-tone discrimination, typical adults exhibit JNDs of around 10-20 cents, and trained listeners can reach 3-4 cents near 1 kHz roughly 0.2% (Bianchi et al., 2016; Mathias, 2010; Zarate et al., 2012; Sun et al., 2017). In SonicBench, we deliberately work in a much easier regime, using 1-semitone (100-cent) contrasts that are trivial for normal-hearing listeners. This coarse, forced-choice design avoids borderline detectability and turns pitch into a foundational probe of whether models encode stable, human-like pitch representations that can in principle support more complex downstream reasoning.

Brightness. Brightness characterizes how sharp or dull a sound is, reflecting how much spectral energy is concentrated at higher versus lower frequencies. This cue is central to timbre and voice perception, equalization and mixing decisions, and everyday source or scene identification (Warren et al., 2005; Välimäki and Reiss, 2016; Town and Bizley, 2013). Psychoacoustic work shows that perceived brightness is strongly linked to spectral centroid and related high-frequency energy measures, and that listeners can reliably order and discriminate changes along this dimension across a wide range of instruments and environmental sounds (Schubert and Wolfe, 2006; Almeida et al., 2017; Saitis and Siedenburtg, 2020). In SonicBench, we do not analytically manipulate the centroid; instead, we draw on community-endorsed corpora with brightness annotations and apply manual screening to ensure that paired items are clearly separable to human listeners (Engel et al., 2017a). This yields a coarse, forced-choice probe of whether models encode a stable notion of spectral balance, rather than merely reacting to extreme or degenerate spectral conditions.

Loudness. Loudness characterizes the perceived intensity of a sound rather than its raw signal level. It is central to speech intelligibility, mix balance, and attentional salience in everyday listening, for example when a listener must follow one talker in noise or react to a sudden, loud alert. Psychophysical studies show that normal-hearing listeners are

highly sensitive to level differences: for broadband or complex stimuli, just-noticeable changes are typically on the order of 0.5–1 dB, and for pure tones the difference limen can drop to roughly 0.2–0.6 dB at higher sensation levels (Florentine et al., 1987; Hall et al., 2008). In SonicBench, we operationalize loudness using integrated LUFS with K-weighting following ITU-R BS.1770 and EBU R128 (itu, 2015; ebu, 2014), and construct contrasts that sit well above these thresholds (e.g., ≥ 2 –3 LU). This coarse, forced-choice design ensures that the task is trivial for human listeners and probes whether models encode a robust representation of perceived loudness, rather than merely skirting the boundary of detectability.

Velocity. Velocity captures how forcefully a sound-producing action is executed, shaping both the sharpness of its onset and its overall energy. In everyday listening, many impact- and contact-driven sounds e.g., knocking, footsteps, percussion, keyboard instruments change systematically with the speed and force of the underlying motion, and listeners rely on these cues to infer effort, material, and even affect (Giordano, 2005; Giordano and McAdams, 2006). In practice, this notion is often operationalized as a control parameter that co-regulates loudness and spectral brightness, for example MIDI velocity (1-127) in digital or computer-controlled instruments, which correlates closely with hammer speed and radiated level on systems such as the Yamaha Disklavier (Goebel and Bresin, 2001). Behavioural studies on impact and performance sounds show that humans are highly sensitive to such variations and can reliably distinguish differences in playing effort and dynamic strength well within this control range (Giordano, 2005; Lemaitre and Heller, 2012). In SonicBench, we instantiate this attribute using recordings with explicit velocity-related control and construct pairs that differ by a substantial margin in excitation strength while keeping pitch and other factors as stable as possible, yielding a coarse, forced-choice probe of whether models encode a robust representation of dynamic impact strength rather than merely reacting to incidental loudness fluctuations.

A.2 Temporal

This dimension captures how acoustic events unfold over time, from the duration of individual segments to the pacing of repeated patterns. Temporal cues are a core low-level property of the signal,

yet humans routinely recruit them for higher-level judgments, using duration and pacing to segment actions, perceive urgency, and distinguish slow, heavy motion from fast, agile motion. For audio understanding models, such temporal understanding is increasingly a foundational requirement, beyond answering explicit timing questions, models in complex settings such as interactive assistants should be able to infer changes in speed, intent, or urgency from patterns like footsteps gradually accelerating or alarms changing rate. We therefore instantiate this dimension with two complementary attributes, Duration and Tempo, which targets sensitivity to absolute interval length and probes how well models track relative pacing across repeated or patterned sounds.

Duration. Duration denotes the physical time span between the onset and offset of an isolated sound event and serves as a basic probe of temporal awareness. Psychophysical studies of auditory interval discrimination in the hundreds-of-milliseconds to seconds range suggest that sensitivity is well approximated by a roughly constant Weber fraction on the order of 10–15%, i.e., a 600 ms reference typically requires a change of about 60–90 ms to be reliably discriminated (Grondin, 1993; Rammsayer, 2010; Rammsayer and Ulrich, 2012). This stable, low-level acuity makes duration a natural baseline for assessing whether models have acquired a human-like temporal scaffold on which more complex timing judgements can build. In SonicBench, we therefore construct duration contrasts that sit well above this regime e.g., differences of 30–50% relative to the base interval, turning the task into a coarse, forced-choice probe of temporal grounding: we test whether models can represent absolute time, segment events, and compare intervals in a robust way, rather than merely responding to fine-grained threshold differences.

Tempo. Tempo captures the rate at which a sequence of events unfolds, i.e., how quickly successive onsets occur over time. It is a property of the temporal pattern rather than an absolute interval: changing tempo alters the spacing between events while preserving their order. Humans routinely recruit tempo cues for higher-level inference—for example, using accelerating footsteps to infer a transition from walking to running, or changes in machine cycles to detect shifts in operating state—so robust tempo perception forms a foundational layer of temporal understanding on which more

complex reasoning about actions and scenes can build. Psychophysical studies indicate that tempo discrimination approximately follows Weber’s law, with just-noticeable differences on the order of a few percent of the base rate roughly 4–8% at moderate tempi (Michon, 1964; Drake and Botte, 1993; Boltz, 1998; ten Hoopen et al., 2004). In SonicBench, we operationalize tempo using rhythmic and musical excerpts with explicit beat annotations and construct pairs whose tempo differences are substantially larger than these human JNDs. This yields a coarse, forced-choice probe of whether models encode stable representations of relative pacing, rather than merely reacting to marginally detectable speed changes.

A.3 Spatial & Environment

This dimension captures how sounds are embedded in space and in their surrounding acoustic environment whether a source is in front or behind, near or far, and in a dry booth or a reverberant hall. For humans, such cues are part of the low-level acoustics of every signal, yet they are constantly recruited for higher-level judgements, orienting toward a talker, judging whether an approaching sound is dangerously close, or inferring that a voice is echoing in a large indoor space. For LALMs and embodied agents, spatial and environmental grounding is therefore a foundational capacity rather than a niche skill, without a basic sense of direction, depth, and room context, models cannot reliably follow situated instructions, coordinate movement, or assess risk in real-world audio scenes. SonicBench instantiates this dimension with three attributes, Direction, Distance, and Reverberation, each framed as a simple binary judgement as Figure 1 shown that isolates one spatial or environmental property while leaving other cues as controlled as possible.

Direction. Direction denotes the perceived azimuth of a sound source where it is located around the listener in the horizontal plane. Humans are highly sensitive to directional cues, especially for frontal sources, achieving minimum audible angles of only a few degrees under favourable conditions (Mills, 1958; Makous and Middlebrooks, 1990), supported by interaural time and level differences together with direction-dependent spectral filtering by the head and pinnae (Blauert and Hearing, 1997). This acuity underlies everyday behaviours such as orienting toward a talker, avoid-

ing approaching vehicles, and monitoring events outside the visual field, and is equally crucial for LALMs deployed in embodied or interactive settings. In SonicBench, we treat direction as a coarse, foundational spatial attribute rather than a fine-grained localisation task. We restrict the horizontal plane to two clearly separated regions: a front sector spanning $\pm 60^\circ$ around 0° and a back sector spanning $\pm 60^\circ$ around 180° . For each stimulus, we sample an azimuth uniformly within one sector and map it to the corresponding front/back label. These large angular separations far exceed human just-noticeable differences, yielding a coarse, forced-choice probe of whether models exhibit stable front-back awareness at all, rather than near-threshold localisation acuity.

Distance. Distance denotes the perceived proximity of a sound source along the source-listener axis—how near or far it seems, rather than its physical distance in metres. Human distance perception draws on a combination of cues, including overall level falloff, spectral changes, and, in typical rooms, the direct-to-reverberant energy ratio (DRR). Psychophysical studies show that listeners can discriminate relatively small changes in source distance at near ranges, and that DRR becomes a dominant cue as distance increases in reverberant environments (Zahorik, 2002; Kolarik et al., 2016). Such proximity judgments are crucial for everyday behaviour and for LALMs in embodied settings, for example when deciding whether a sound source is close enough to warrant attention or action. In SonicBench, we therefore treat distance as a coarse, foundational attribute rather than a fine-grained ranging task. We partition stimuli into two broad bands, near versus far, separated by substantial changes in nominal source distance and DRR, and label them accordingly. These large separations lie well above typical human discrimination thresholds and yield a simple forced-choice probe of whether models encode any stable sense of auditory proximity.

Reverberation. Reverberation captures whether a sound is perceived as occurring in a reflective environment or in an acoustically dry, near-anechoic setting. Psychophysically, reverberation arises from dense patterns of early reflections and late decay, and listeners are highly sensitive to these cues: even modest changes in reverberation time, direct-to-reverberant ratio (DRR), or early reflection structure can alter the perceived room size,

source–room distance, and overall sense of envelopment (Zahorik et al., 2005; Kuttruff, 2016). Such environmental impressions are critical for everyday reasoning about whether a sound is indoors or outdoors, in a small room or a large hall, or mediated by a strongly reflective space. In SonicBench, we therefore treat reverberation as a coarse, foundational environmental attribute rather than a fine-grained RT60 or DRR estimation task. Each item is constructed either as a dry signal with negligible room contribution or as a reverberant version obtained by convolving the same source with a realistic room impulse response, while keeping source content, level, and direction as stable as possible. This binary, far-above-threshold contrast turns the task into a robust forced-choice probe of whether models can reliably distinguish “in-room” versus “dry” conditions, a prerequisite for downstream behaviours such as inferring indoor/outdoor context, room scale, and whether a sound is likely mediated by an enclosed space.

A.4 Timbre

In contrast to the largely low-level physical cues above, the timbre dimension targets how spectral and temporal structure are integrated into mid-level representations of “what kind of sound this is.” Timbre allows listeners to tell a flute from a violin at the same pitch and loudness, or to distinguish crackling fire from rustling leaves, by jointly using spectral shape, temporal envelope, and fluctuations. This ability underpins everyday source and material recognition and acts as a bridge between raw acoustics and semantic understanding where an agent that cannot reliably infer “what is making this sound” or “what kind of background this is” will struggle to interpret scenes or follow natural, audio-grounded instructions. In SonicBench, we therefore instantiate this dimension with two complementary attributes, Timbre, which probes source-centred identity under controlled pitch and loudness, and Texture, which targets the statistical structure of dense, background-like sound fields such as rain, fire, or crowd noise.

Timbre. Timbre is the attribute of auditory sensation that allows listeners to judge two sounds as different even when their pitch and loudness are matched. Perceptually, it reflects a mid-level integration of spectral shape, temporal envelope, and fine-structure cues, and underlies our ability to recognise “what is sounding” (e.g., instrument

type, material, or mode of excitation) from very short excerpts. Human listeners are highly proficient at this kind of judgement: even brief, isolated notes often suffice to distinguish familiar instruments or everyday sound sources. In SonicBench, we operationalise timbre as source-centred discrimination under controlled conditions. We draw on existing corpora with reliable instrument labels, normalise pitch and loudness within each pair, and construct contrasts where the primary difference lies in timbral character rather than in fundamental frequency or overall level. Our goal is not to introduce yet another general-purpose instrument-classification benchmark, but to provide a pragmatic, well-controlled proxy for timbre perception within our broader taxonomy of physical attributes. This coarse, forced-choice setup probes whether LALMs encode robust representations of source identity from timbral cues alone, while remaining flexible enough to be extended beyond the specific instrument categories used here.

Texture. Texture refers to aggregate, noise-like sound patterns produced by the superposition of many similar acoustic events, such as rainfall, crackling fire, or dense insect choruses (McDermott and Simoncelli, 2011). Unlike timbre, which is tied to the identity of a single, salient source, auditory texture perception is inherently statistical: listeners do not track individual droplets or clicks, but summarise ensemble properties over time e.g., density, spectral spread, modulation patterns and use these summaries to recognise and categorise backgrounds and environments in everyday listening. This mid-level representation provides a bridge from raw acoustics to semantic scene understanding, supporting judgements such as whether a setting sounds “rainy”, “crowded”, or “mechanical” even when no single event is isolated. In SonicBench, we operationalise texture as texture-centred environment discrimination. We draw on existing corpora with reliable labels for prototypical textures e.g., rain, fire, crowd noise and construct pairs in which the dominant difference lies in their texture statistics, while overall level and other basic factors are normalised as far as practicable. Our goal is not to build a full-scale environmental sound classification benchmark, but to provide a coarse, forced-choice probe of whether LALMs encode stable statistical representations of background texture that can support higher-level scene reasoning.

A.5 Scene-Level

This dimension targets properties of the entire auditory scene rather than individual sources or isolated events. Whereas the previous dimensions probe local cues such as spectrum, timing, spatial layout, and timbre, scene-level perception asks how many distinct things are happening and how they are organised over time. In human hearing, this corresponds to classic auditory scene analysis (Bregman, 1994), grouping sounds into events and streams, then summarising them at an abstract level e.g., “three knocks before the door opens” or “several dogs barking at once”. For LALMs and embodied agents, such holistic structure is a foundational layer above low-level detection, supporting decisions about complexity, crowding, and how many entities may require attention or action. In SonicBench, we instantiate this dimension with one key attribute, Counting which probes whether models can move beyond simple “is it there?” detection to track how many target events occur in a soundscape.

Counting. Counting assesses the ability to enumerate repeated sound events of a given type within a short clip e.g., door knocks, dog barks. This goes beyond simple detection: a system must segment the waveform into discrete events, maintain them in a working representation, and compare their total number. Human studies suggest that auditory numerosity follows a pattern similar to visual “subitizing”: small numerosities around 1–3 or 4 events can be judged rapidly and accurately, while performance and response times degrade as the count increases and attention and working memory become limiting resources (Cowan, 2001; Anobile et al., 2021). Guided by these findings, SonicBench uses stimuli containing between 1 and 6 target events with controlled event type and spacing, and formulates the task as a coarse multiple-choice judgement over candidate counts. This design keeps the contrasts well within the range that is straightforward for human listeners, while providing a scene-level probe of whether models can move beyond binary detection to robustly track *how many* distinct events occur in a soundscape.

B Toolbox Details

B.1 Mode Selection

The SonicBench Toolbox supports both modes of use: it accepts short raw audio clips as the pri-

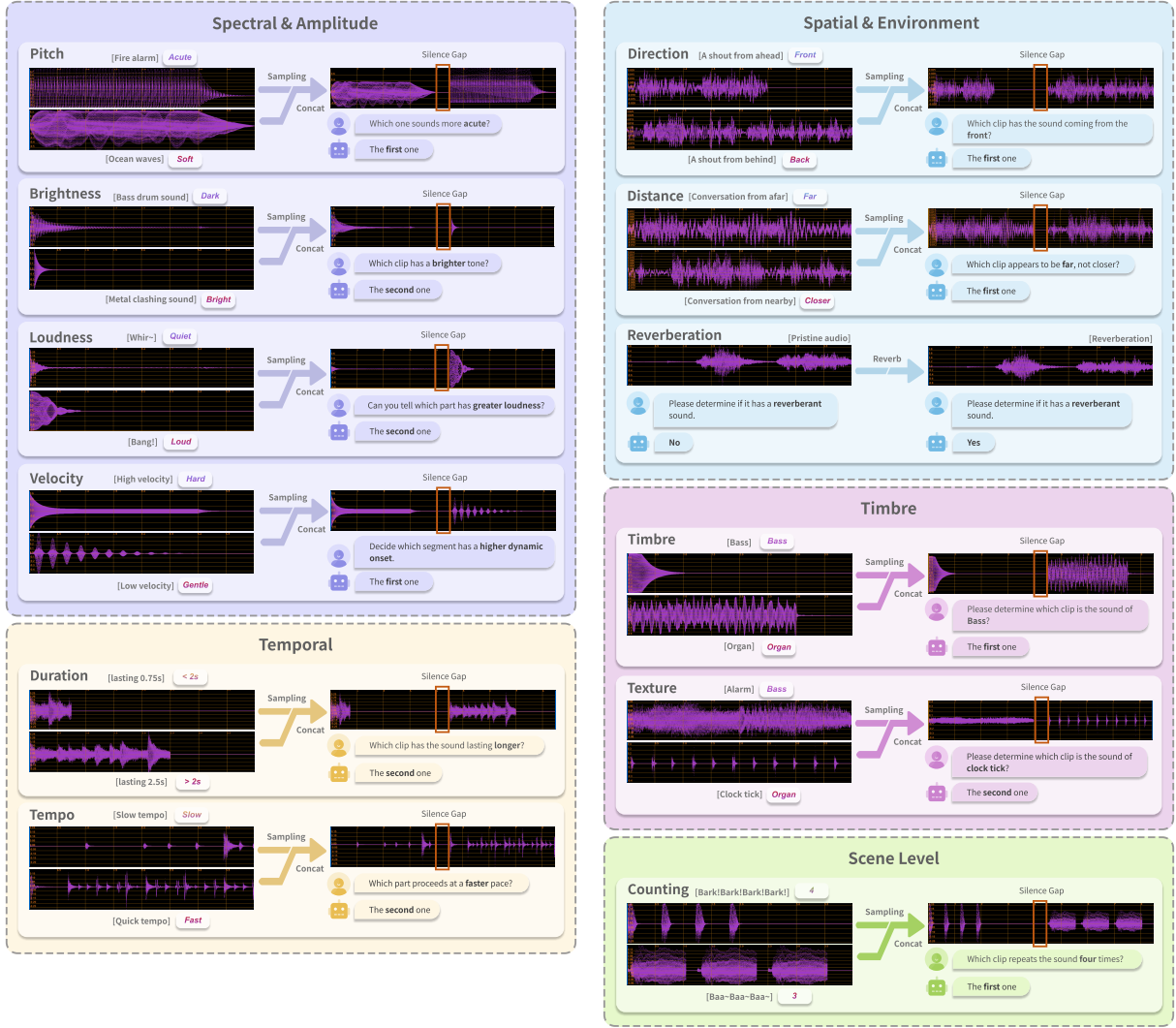


Figure 8: Overview of large scale sampling modes from toolbox used in the dataset construction.

many acoustic input for user customized generation, and it also enables large scale sampling from extensive raw audio collections under user specified constraints for automatic batch stimulus creation. Both modes share a unified input specification and automatic preprocessing pipeline that ensure consistency, experimental control, and reproducibility. Unless otherwise noted, input clips must be mono, 16-bit PCM, and between 0.5 and 5.0 seconds in duration. The clips should contain a single, stable sound event such as an instrument note or isolated vocal token to serve as a reliable reference. Upon ingestion, the toolbox automatically resamples, gain normalizes, and trims or pads the signal to a 4.0 second, 48 kHz reference waveform. This removes the need for manual editing in a DAW and guarantees uniform downstream processing.

User Customized Mode. In this mode, we provides reference clips together with a target config-

uration, such as the desired pitch shift, duration, or loudness level. The toolbox then applies deterministic rule based signal processing recipes that modify only the specified attribute while holding all other acoustic dimensions constant. Users are encouraged to provide clean and monophonic reference recordings to ensure that intended single attribute contrasts remain unconfounded.

Large Scale Sampling Mode. (Used in This Work) In this mode, the user supplies a task specification instead of reference audio. The minimal specification includes the task type, such as comparison or recognition, the target attribute or attributes selected from the twelve item taxonomy (pitch, brightness, loudness, velocity, duration, tempo, direction, distance, reverberation, timbre, texture, counting), value ranges or sampling constraints, and any language or format requirements, as shown in Figure 8. The toolbox samples attribute values

within the defined space and synthesizes the corresponding audio automatically, producing example pairs, evaluation prompts, and gold annotations suitable for large experimental sets.

B.2 Rule-Based Attribute Control

Fine-grained Attribute Control. In user customized mode, the toolbox supports fine-grained adjustment of all twelve target attributes while keeping other acoustic dimensions constant. Pitch is shifted using PSOLA or phase vocoder, brightness is adjusted via spectral envelope transformations, and loudness is controlled using EBU R128 calibrated gain. Velocity is modulated through amplitude envelope control, duration and tempo are modified via WSOLA or time stretching, and directional cues are changed with HRTF or panning. Distance is adjusted by controlling the direct-to-early-reflection ratio with loudness compensation, reverberation is altered via convolution and artificial reverb, and timbre is shaped through filtering. Texture is modified using noise enhancement, modulation, or granular processing, and counting is controlled by duplicating, trimming, or rhythmically arranging discrete events. These rule-based operations ensure that each generated stimulus isolates the intended attribute for precise perceptual comparisons.

Attribute-controlled Sampling. For six continuous and categorical attributes including pitch, loudness, velocity, tempo, duration, and distance, we first compute feature statistics for all candidate samples and discretize them into bins using perceptually motivated thresholds. Specifically, pitch is divided using a boundary at 65, loudness is thresholded at -15 LUFS, velocity at 75, duration at 2.4s, and tempo at 100 BPM. Stratified sampling is then performed based on these bins. Candidate clips that meet the specified criteria are selected, and samples are drawn across multiple bins to ensure perceptible differences while preserving the pre-sampling bin distribution. This procedure eliminates potential sampling biases and ensures that non-target attributes remain strictly controlled, avoiding confounding effects in downstream evaluation.

For brightness, candidate clips are filtered from a large pool to include only clearly perceptible bright or dark samples, covering a wide variety of sound types such as musical instruments, environmental sounds, human speech, and animal vocalizations. Duration-controlled sampling extracts valid audio

segments from diverse sources including music, animal calls, daily conversations, and speeches, which are then zero-padded to create uniform 4-second clips.

For reverberation, clean audio is convolved with different room impulse responses to produce varying reverb levels. During evaluation, all other attributes and content are held constant, with only reverberation varying. For counting, the number of sound events is controlled by replicating audio segments to achieve between 1 and 7 occurrences, such as repeated animal calls or impact sounds, while keeping all other attributes fixed.

Texture control involves automated sampling from audio with distinct spectral or temporal roughness patterns, producing a rich set of clips that span smooth, granular, and noisy textures, suitable for perceptual discrimination tasks. For directional attributes, we define a broad front sector as within $\pm 60^\circ$ from the center and a broad rear sector similarly, ensuring that source positions fall within these regions while allowing controlled variation in angular difficulty.

B.3 Task Construction

For recognition tasks, the toolbox generates a single transformed audio clip according to the specified target attribute and automatically produces accompanying textual instructions, categorical labels, and ground truth answers. This fully automated process allows models to accurately identify the direction, magnitude, or category of the attribute change without manual intervention. For comparison tasks, the system automatically constructs paired datasets consisting of the original and transformed audio clips and generates explicit comparative prompts, such as “Which clip is brighter,” “Which clip is longer,” or “Which clip has a higher pitch.” Corresponding ground truth answers are also provided, enabling straightforward, standardized, and reproducible evaluation of model performance. The design supports flexible experimentation across all twelve defined attributes and can be scaled to generate large volumes of controlled stimuli for both recognition and comparison scenarios.

B.4 Output

Each run of the toolbox produces the following outputs. First, audio files including both reference and transformed samples, with filenames explicitly indicating the attribute changes. Second, a task annotation file in JSON format that records the

task type, target attribute, audio file paths, textual instructions, and ground truth answers. Third, a batch-level index in CSV or JSONL format to facilitate downstream evaluation, large-scale processing, and dataset release. These outputs ensure full traceability, reproducibility, and ease of integration into automated evaluation pipelines.

B.5 Reproducibility

All signal processing operations in the toolbox are executed deterministically, ensuring that given the same input audio and configuration, the generated outputs are bitwise reproducible. This guarantees that SonicBench produces consistent stimuli across repeated runs, enabling fair and reliable evaluation. Such reproducibility provides a solid foundation for comparing different models, languages, or dataset scales, and ensures that experimental conclusions remain consistent regardless of the computational environment or deployment scenario.

C Annotator Qualifications

Although SonicBench builds upon curated open-source corpora, the critical phases of data filtration, concatenation, and final perceptual verification required expert human judgment to ensure physical correctness. To maintain the highest standards of quality control, our team was composed entirely of researchers with deep academic backgrounds. Every team member held at least a bachelor’s degree, with over a half currently pursuing doctoral studies in relevant fields. Crucially, all personnel were screened for normal hearing function, ensuring they possessed the auditory acuity necessary to validate JNDs and confirm the perceptual salience of attribute contrasts across all samples.

D Details of Data Sources

In this section, we provide detailed specifications of the data sources utilized to construct SonicBench. The selection logic prioritizes datasets that offer verifiable ground-truth labels and distinct signal properties relevant to the target dimension. Table 5 provides a systematic mapping between the evaluation dimensions and their corresponding source datasets. The specific characteristics of each dataset are detailed below.

NSynth. (Engel et al., 2017a) It’s a large-scale, high-quality dataset for neural audio synthesis of musical notes, consisting of approximately 306043 four-second monophonic audio snippets (sampled

| Dimensions | Attributes | Data Source |
|-----------------------|---------------|------------------------------------|
| Spectral & Amplitude | Pitch | NSynth (Engel et al., 2017a) |
| | Brightness | |
| | Loudness | |
| | Velocity | |
| Temporal | Duration | TAU (Adavanne et al., 2019) |
| | Tempo | Groove MIDI (Gillick et al., 2019) |
| Spatial & Environment | Direction | LOCATA (Evers et al., 2020) |
| | Distance | RealMAN (Yang et al., 2024a) |
| | Reverberation | SRIRS (Amengual Garí et al., 2020) |
| Timbre | Timbre | NSynth (Engel et al., 2017a) |
| | Texture | ESC-50 (Piczak, 2015) |
| Scene Level | Counting | PicoAudio (Xie et al., 2024) |

Table 5: Overview of Data Source.

at 16kHz) from around 1000 musical instruments. Each note is annotated with pitch, quality, velocity, brightness, etc.

TAU. (Adavanne et al., 2019) This dataset is a multi-room reverberant dataset designed for Sound Event Localization and Detection (SELD) in the DCASE 2019 challenge. It includes two subsets, each with a 400-recording development set and a 100-recording evaluation set. It is synthesized by convolving isolated sound events from DCASE 2016 Task 2 with real impulse responses from 5 indoor environments, adding ambient noise from these environments. This dataset is labeled with azimuth, elevation angle, start and end time.

Groove MIDI. (Gillick et al., 2019) It’s a large-scale dataset for expressive drum performance research, containing 13.5 hours of recordings of drummers playing electronic drum kits. It includes detailed metadata (drummer IDs, musical style, tempo), and is an order of magnitude larger than the largest previously publicly available comparable dataset for drum performance modeling.

LOCATA. (Evers et al., 2020) This dataset is an open-access corpus tailored for benchmarking acoustic source localization and tracking algorithms, containing multichannel audio recordings from four distinct microphone arrays. It is fully annotated with critical information including ground-truth positions/orientations of sources and sensors, and hand-labelled voice activity periods.

RealMAN. (Yang et al., 2024a) RealMAN is a real-recorded and annotated microphone array dataset designed for dynamic speech enhancement and source localization, which uses a 32-channel high-fidelity microphone array to record speech across 32 scenes and background noise across. It

provides key annotations such as direct-path target clean speech and speaker location.

SRIRS. (Amengual Garí et al., 2020) It’s created by Detmold University of Music, comprises around 600 multichannel Spatial Room Impulse Responses (SRIRs) captured in three spaces. It includes unique setups such as artificial reverberation in Detmold Konzerthaus, Wave Field Synthesis (WFS) focused sources, and stage measurements with music stands.

ESC-50. (Piczak, 2015) The ESC Dataset, created for environmental sound classification research. It derives from Freesound and standardized to 44.1 kHz, single-channel Ogg Vorbis format. It also provides human classification accuracy benchmarks and baseline machine learning results, along with replication code via a Jupyter notebook, supporting open research in auditory recognition.

PicoAudio. (Xie et al., 2024) The dataset used in the PicoAudio project is constructed by crawling audio clips from Freesound using sound event keywords, then segmenting them with a text-to-audio grounding model and filtering via the LAION-CLAP model. Each sample is paired with timestamp captions and frequency captions to support training and evaluating temporally controllable audio generation models.

E Data Selection and Filtering

E.1 Data Selection

In constructing SonicBench, we first perform a systematic data selection process over existing public audio datasets to identify samples that are suitable for evaluating fundamental physical auditory perception. The primary selection criterion is the definability and controllability of target physical attributes. Although the audio samples are not directly generated by us, we only retain those for which key physical properties, can be reliably estimated or are annotated by the original datasets. Moreover, the values of these attributes must be well-defined and stable, enabling consistent formulation of both absolute (recognition) and relative (comparison) judgment tasks. For continuous attributes, we prioritize samples whose values exhibit clear distributions and smooth variations that do not rely on semantic interpretation. For discrete or graded attributes, category boundaries must be grounded in consistent physical criteria rather than

subjective or semantic definitions, ensuring that all labels retain clear physical meaning.

Building upon this, the data selection process further follows the principles of single-attribute dominance and task compatibility. Each selected sample is used to evaluate only one target physical attribute, minimizing confounding effects from simultaneous variations in multiple cues. For comparison tasks, samples are organized into pairs such that only the target attribute differs systematically, while all other physical attributes are held constant or remain statistically symmetric. In addition, we retain only samples that naturally support both absolute (recognition) and relative (comparison) paradigms, allowing the same physical attribute to be evaluated consistently under single-stimulus identification and pairwise discrimination settings.

E.2 Data Filtering

After completing data selection, we apply a systematic data filtering process to the retained samples as a form of secondary validation, removing instances that may still compromise the effectiveness of the evaluation. We first filter the data from the perspective of perceptual validity, with particular emphasis on excluding audio samples in which the target physical attribute is insufficiently salient or exhibits unstable values. For continuous physical attributes, we remove samples whose attribute differences are too small, lie close to perceptual boundary conditions, or show high uncertainty in the time or frequency domains. For spatial, environmental, and timbral attributes, we exclude cases where the available acoustic cues are insufficient to induce stable and reliable perceptual differences. This step ensures that, for all retained samples, the target physical attribute is clearly and discriminably expressed at the signal level, preventing model errors from being erroneously attributed to ambiguities in the input audio itself.

The data filtering process further enforces strict control over confounding factors, label consistency, and task structure. For comparison tasks, we remove paired samples in which systematic differences remain in non-target attributes, ensuring that model decisions cannot rely on spurious cues. For absolute judgment tasks, we discard samples whose labels do not align cleanly with the corresponding attribute value ranges or exhibit category overlap, thereby preserving a one-to-one correspondence between physical parameters and labels. In addition, we filter out audio samples containing salient

semantic cues, abnormal energy distributions, clipping, or other signal artifacts that could enable shortcut solutions.

F Review & Double Check

As illustrated in the “Review & Double Check” stage of Figure 2, we implemented a rigorous, multi-stage verification pipeline to ensure that every sample in SonicBench adheres to both physical correctness and perceptual validity. This process involves three distinct checkpoints, executed by independent domain experts to enforce role separation.

1. Correction (Metadata & Logic Verification).

The first line of defense involves a manual review of the generated metadata and Question-Answer (QA) pairs. In this phase, an independent annotator (distinct from the author) verifies the consistency between the audio filename parameters and the JSON labels. Common errors targeted in this stage include label mismatches (e.g., a file generated with high pitch labeled as low) or grammatical inconsistencies in the instruction templates. Samples containing correctable errors are routed to revision, while those with fundamental logical flaws are flagged for rejection.

2. Significant Difference Check (Perceptual Validation).

Passing samples proceed to the most critical phase: the Significant Difference Check. Here, expert annotators listen to the audio to confirm that the target attribute manipulation is perceptually salient. For comparison tasks, this step rigorously enforces the JND constraint. Annotators must verify that the contrast between the two clips is immediately distinguishable to a healthy human ear without ambiguity. For recognition tasks, annotators verify that the attribute value (e.g., “High Brightness”) clearly aligns with the auditory sensation. As shown in the pipeline diagram, any sample failing this perceptual threshold, where the difference is too subtle or masked by other acoustic features, is immediately routed to the “Rejected Samples” pool to prevent “guessing games” in evaluation.

3. Quality Check (Signal Integrity). The final gate is a technical Quality Check. This step focuses on the acoustic fidelity of the file. Annotators inspect the waveform for generation artifacts, such as:

- **Boundary Artifacts:** Unnatural clicks or pops at the splicing points (0.5s silence intervals) in comparison pairs.
- **Digital Distortion:** Clipping or excessive synthetic noise introduced during the rendering process.
- **Format Compliance:** Ensuring duration (4s/8.5s) and sampling rate consistency.

Only samples that pass this final inspection move to the final benchmark

Rejection Policy. As depicted by the red dashed paths in Figure 2, a “Fail” at any of the three stages triggers a rejection mechanism. To ensure high quality, we adopted a strict policy: while minor metadata issues allowed for one round of correction, any failure regarding perceptual salience (Stage 2) or irreversible signal degradation (Stage 3) resulted in permanent discarding. After this process, we retained 2,400 high-quality questions in the final version. We provide our benchmark statistics in Table 2.

G Details of Task Instruction

| | |
|---|---|
| <p>Pitch</p> <p><i>(Recognition Task)</i></p> <p>Question: "Listen to the audio. Pitch is high if above 65, otherwise low. Is the sound high-pitched or low-pitched? Only reply 'A' if high-pitched, 'B' if low-pitched. Only answer 'A' or 'B'. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.high-pitched B.low-pitched</p> <p><i>(Comparison Task)</i></p> <p>Question: "Two auditory segments are concatenated with 0.5 seconds of silence. Decide which segment has a higher pitch. Only reply 'A' first, 'B' second. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.first B.second</p> | <p>Brightness</p> <p><i>(Recognition Task)</i></p> <p>Question: "Is the audio's spectral energy skewed toward higher frequencies (bright) or lower (dark)? Enter 'A' for brightness, 'B' for dark — nothing else. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.brightness B.dark</p> <p><i>(Comparison Task)</i></p> <p>Question: "The audio includes two segments with a 0.5-second silent interval. Which is brighter? Only answer letter 'A' (refers to the first clip) or 'B' (refers to the second clip). Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.the first clip B.the second clip</p> |
| <p>Loudness</p> <p><i>(Recognition Task)</i></p> <p>Question: "How would you describe the volume of this audio: loud or soft? Choose 'A' if loud, 'B' if soft. Consider a sound loud if its KUFS value is greater than -15, otherwise soft. Only answer 'A' or 'B'. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.loud B.soft</p> <p><i>(Comparison Task)</i></p> <p>Question: "A dual-segment audio is presented with 0.5 seconds of silence in between. Which segment has the higher volume? Use 'A' for the first, 'B' for the second. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.first B.second</p> | <p>Velocity</p> <p><i>(Recognition Task)</i></p> <p>Question: "A single audio clip is presented. Classify the velocity as fast if it is greater than 75; otherwise, classify it as slow. Determine if the note was struck hard or softly. Answer 'A' for high, 'B' for low. Only answer 'A' or 'B'. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.high B.low</p> <p><i>(Comparison Task)</i></p> <p>Question: "Listen carefully: two clips are separated by 0.5 seconds of silence. Compare the perceived intensity of the two clips. Answer 'A' if the first is stronger, 'B' if the second. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.first B.second</p> |

Figure 9: Examples of Spectral & Amplitude Dimension. Shown are representative samples for Pitch, Brightness, Loudness, and Velocity attributes across both Recognition and Comparison tasks.

| | |
|--|---|
| <p>Duration</p> <p><i>(Recognition Task)</i></p> <p>Question: "Can you determine if the sound is longer than 2.4 seconds? Reply with 'A' if it is longer than 2.4 seconds, 'B' if it is equal to or shorter than 2.4 seconds. Only answer 'A' or 'B'. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.longer than 2.4 seconds B.shorter than 2.4 seconds</p> <p><i>(Comparison Task)</i></p> <p>Question: "Two audio clips are joined by a short silence. Do you know which clip is longer? Enter 'A' for longer first, 'B' for longer second. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.longer first B.longer second</p> | <p>Tempo</p> <p><i>(Recognition Task)</i></p> <p>Question: "Evaluate the speed of the sound: is it fast or slow? Use 'A' to mean fast tempo, 'B' to mean slow tempo. Consider the tempo fast if it is greater than 100; otherwise, consider it slow. Only answer 'A' or 'B'. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.greater than 100 B.less than 100</p> <p><i>(Comparison Task)</i></p> <p>Question: "Two auditory recordings are given, separated by a small pause. Which part is performed slower? Only answer 'A' for the first, 'B' for the second. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.first B.second</p> |
|--|---|

Figure 10: Examples of Temporal Dimension. Shown are representative samples for Duration and Tempo attributes across both Recognition and Comparison tasks.

| | |
|--|---|
| <p>Reverberation</p> <p><i>(Recognition Task)</i></p> <p>Question: "An audio sample is provided. Judge whether it contains reverberation. Say 'A' if reverberant, 'B' if not. Only answer 'A' or 'B'. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.reverberant B.not reverberant</p> <p><i>(Comparison Task)</i></p> <p>Question: "You are provided with two clips and a 0.5-second silent interval. Assess which segment appears to be in a more reverberant space. Output only 'A' (first) or 'B' (second). Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.first B.second</p> | <p>Distance</p> <p><i>(Recognition Task)</i></p> <p>Question: "You are given an audio clip. Please determine if the sound source is without 3 meters. Enter 'A' if the sound is near (within 3m), 'B' if far. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.within 3m B.without 3m</p> <p><i>(Comparison Task)</i></p> <p>Question: "Which clip sounds like it's not close but distant? Only answer letter 'A' (refers to the first clip) or 'B' (refers to the second clip). Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.first clip B.second clip</p> |
| <p>Direction</p> <p><i>(Recognition Task)</i></p> <p>Question: "Is the sound more front or back? Only answer letter 'A' if it is from the front, 'B' if it is from the back. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.from the front B.from the back</p> <p><i>(Comparison Task)</i></p> <p>Question: "Which segment seems to be located in the front of the stereo field? Only answer letter 'A' (refers to the first clip) or 'B' (refers to the second clip). Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.first clip B.second clip</p> | <p>Counting</p> <p><i>(Recognition Task)</i></p> <p>Question: "How many times is the sound produced? Only answer letter 'A' for 4, 'B' for 1. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.4 B.1</p> <p><i>(Comparison Task)</i></p> <p>Question: "You are given a concatenated audio of two clips with 0.5 seconds of silence in between. Please determine which clip is the sound of toilet_flush? Only answer letter 'A' (refers to the first clip) or 'B' (refers to the second clip). Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.first clip B.second clip</p> |

Figure 11: Examples of Spatial & Environment and Scene Level Dimensions. Shown are representative samples for Reverberation, Distance, Direction, and Counting attributes across both Recognition and Comparison tasks.

| | |
|---|--|
| <p>Texture</p> <p><i>(Recognition Task)</i></p> <p>Question: "Can you identify the sound as A. clapping or B. drinking sipping? Only answer 'A' or 'B'. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.clapping B.drinking sipping</p> <p><i>(Comparison Task)</i></p> <p>Question: "You are given a concatenated audio of two clips with 0.5 seconds of silence in between. Please determine which clip is the sound of toilet_flush? Only answer letter 'A' (refers to the first clip) or 'B' (refers to the second clip). Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.first clip B.second clip</p> | <p>Timbre</p> <p><i>(Recognition Task)</i></p> <p>Question: "Does it sound like A. reed or B. guitar? Only answer 'A' or 'B'. Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.reed B.guitar</p> <p><i>(Comparison Task)</i></p> <p>Question: "Consider the following audio clip. Please determine which clip is the sound of keyboard? Only answer letter 'A' (refers to the first clip) or 'B' (refers to the second clip). Do not add any explanation, punctuation, or extra text. <audio>"</p> <p>Choices & Answer: A.first clip B.second clip</p> |
|---|--|

Figure 12: Examples of Timbre Dimension. Shown are representative samples for Timbre and Texture attributes across both Recognition and Comparison tasks.

H Benchmarking Candidates

In this section, we provide a comprehensive overview of the models evaluated in SonicBench. Our assessment covers a diverse spectrum of architectures, categorized into LALMs, LARMs, and OLMs, encompassing both open-source community checkpoints and state-of-the-art proprietary services. Detailed specifications for each candidate, including specific model names, weight versions, access URLs, and parameter counts, are systematically tabulated in Table 6.

H.1 Large Audio Language Models

Kimi-audio. (KimiTeam et al., 2025) Developed by Moonshot AI. It’s an audio foundation model featuring a 12.5Hz audio tokenizer (combining discrete semantic tokens and continuous acoustic vectors), an LLM-based core, and a chunk-wise streaming detokenizer via flow matching. It achieves state-of-the-art performance across various benchmarks.

Qwen2-audio. (Chu et al., 2024) From Alibaba Group’s Qwen Team, it combines a Whisper-large-v3 audio encoder and the Qwen-7B LLM, capable of processing diverse audio inputs and supporting two seamless interaction modes (Voice Chat and Audio Analysis) without system prompts for switching. It outperforms SOTAs on AIR-Bench’s audio-centric instruction-following tasks.

VITA-audio. (Long et al., 2025) It’s co-developed by multiple institutions, an end-to-end large audio model featuring lightweight Multiple Cross-modal Token Prediction (MCTP) modules, which enable generating multiple audio tokens in one model forward pass to achieve zero audio token delay and 3-5 times inference speedup. It outperforms open-source models of similar size on ASR, TTS, and SQA tasks.

Audio Flamingo2. (Ghosh et al., 2025) Co-developed by NVIDIA and UMD, it’s a unified audio-language model designed for universal audio understanding and generation. It adopts a cross-modal alignment framework that integrates a dual-stream audio encoder, a pre-trained LLM, and a lightweight adapter for efficient feature fusion. It supports a 16K context window for audio files up to 20 minutes, enabling multi-turn audio-text interaction. It outperforms baselines on 12 out of 15 audio-language benchmarks, while maintaining

30% higher inference efficiency than other models via its streamlined adapter design.

Voxtral. (Liu et al., 2025a) Developed by Mistral AI, Voxtral features a 32K context window to handle audio files up to 40 minutes and long multi-turn conversations. It delivers SOTA performance in speech transcription, translation, and understanding, surpassing some closed-source models like GPT-4o mini and Gemini 2.5 Flash in specific tasks while maintaining strong text capabilities.

Baichuan-audio. (Li et al., 2025b) From Baichuan Inc., it’s an end-to-end LALM equipped with an 8-layer RVQ Baichuan-Audio-Tokenizer to retain semantic and acoustic information, and an independent audio head for audio token processing. It supports high-quality real-time bilingual (Chinese and English) speech interaction, outperforms peer models in tasks like ASR, TTS, and audio QA.

MiDashengLM. (Dinkel et al., 2025) Developed by Xiaomi Inc., it integrates the open-source Dasheng audio encoder to unify speech, sound, and music into holistic textual representations. Exclusively trained on publicly available datasets, it delivers up to 4 times faster time-to-first-token (TTFT) and 20 times higher throughput than comparable models, while outperforming baselines like Qwen2.5-Omni-7B and Kimi-Audio-Instruct-7B across audio captioning, QA, and paralinguistic classification tasks.

Llama-Omni. (Fang et al., 2025a) This model integrates a pretrained Whisper-large-v3 speech encoder, a trainable speech adaptor, Llama-3.1-8B-Instruct LLM, and a streaming non-autoregressive speech decoder, enabling direct simultaneous generation of text and speech responses from speech instructions without transcription. It achieves a response latency as low as 226ms, efficiently trained on the InstructS2S-200K dataset, outperforms baselines like SpeechGPT in both content and style for speech interaction.

Llama-Omni2. (Fang et al., 2025b) The successor to LLaMA-Omni, it’s built on Qwen2.5 LLM, integrating Whisper’s speech encoder and an autoregressive streaming speech decoder (including a text-to-speech language model and a causal flow matching model) to enable high-quality real-time speech interaction. It outperforms SOTAs in spoken question answering and speech instruction following tasks.

| Model Names | Weight Versions | URLs | #Params |
|--|--------------------------------------|----------------------------|-----------|
| Large Audio Language Models (LALMs) | | | |
| Kimi-Audio-Instruct (KimiTeam et al., 2025) | moonshotai/Kimi-Audio-7B-Instruct | model card | 7B |
| Qwen2-Audio-Instruct (Chu et al., 2024) | Qwen/Qwen2-Audio-7B-Instruct | model card | 7B |
| VITA-Audio-Plus-Vanilla (Long et al., 2025) | VITA-MLLM/VITA-Audio-Plus-Vanilla | model card | 8B |
| Audio Flamingo 2 (Ghosh et al., 2025) | nvidia/audio-flamingo-2 | model card | 3B |
| Audio Flamingo 3 (Goel et al., 2025) | nvidia/audio-flamingo-3 | model card | 7B |
| Voxtral-Mini (Liu et al., 2025a) | mistralai/Voxtral-Mini-3B-2507 | model card | 3B |
| Baichuan-Audio-Instruct (Li et al., 2025b) | baichuan-inc/Baichuan-Audio-Instruct | model card | 7B |
| MiDashengLM (Dinkel et al., 2025) | mispeech/midashenglm-7b-0804-fp32 | model card | 7B |
| Llama-Omni (Fang et al., 2025a) | ICTNLP/Llama-3.1-8B-Omni | model card | 8B |
| Llama-Omni2 (Fang et al., 2025b) | ICTNLP/LLaMA-Omni2-7B | model card | 7B |
| GLM-4-Voice (Zeng et al., 2024) | zai-org/glm-4-voice-9b | model card | 9B |
| SALMONN (Tang et al., 2024) | tsinghua-ee/SALMONN-7B | model card | 7B |
| MU-LLaMA (Liu et al., 2023) | mu-llama/MU-LLaMA | model card | 7B |
| BAT (Zheng et al., 2025b) | - | model card | 7B |
| R1-AQA (Li et al., 2025a) | mispeech/r1-aqa | model card | 7B |
| Step-Audio 2 mini (Wu et al., 2025) | stepfun-ai/Step-Audio-2-mini | model card | 7B |
| MiMo-Audio-Instruct (Xiaomi, 2025) | XiaomiMiMo/MiMo-Audio-7B-Instruct | model card | 7B |
| GPT-4o-Audio (OpenAI, 2024) | gpt-4o-audio-preview-2025-06-03 | model card | - |
| Large Audio Reasoning Models (LARMs) | | | |
| Mellow (Deshmukh et al., 2025) | soham97/mellow | model card | 167M |
| Audio-Reasoner (Xie et al., 2025) | zhifeixie/Audio-Reasoner | model card | 7B |
| GAMA (Ghosh et al., 2024) | sonalkum/GAMA | model card | 7B |
| R1-AQA (think mode) (Li et al., 2025a) | mispeech/r1-aqa | model card | 7B |
| Audio Flamingo 2 Sound-CoT (Kong et al., 2025) | nvidia/audio-flamingo-2-SoundCoT | model card | 3B |
| Audio Flamingo 3 (think mode) | nvidia/audio-flamingo-3-hf | model card | 7B |
| Step-Audio 2 mini Think | stepfun-ai/Step-Audio-2-mini-Think | model card | 7B |
| MiMo-Audio (think mode) | XiaomiMiMo/MiMo-Audio-7B-Instruct | model card | 7B |
| Step-Audio-R1 | stepfun-ai/Step-Audio-R1 | model card | 33B |
| Omni Language Models (OLMs) | | | |
| Qwen2.5-Omni (Xu et al., 2025a) | Qwen/Qwen2.5-Omni-7B | model card | 7B |
| Baichuan-Omni-1.5 (Li et al., 2025c) | baichuan-inc/Baichuan-Omni-1d5 | model card | 7B |
| VITA-1.5 (Fu et al., 2025) | VITA-MLLM/VITA-1.5 | model card | 7B |
| Ola (Liu et al., 2025c) | THUdyh/Ola-7b | model card | 7B |
| HumanOmni (Zhao et al., 2025) | StarJiaxing/HumanOmni-7B | model card | 7B |
| OpenOmni (Luo et al., 2025) | Tongyi-ConvAI/OpenOmni | model card | 7B |
| Qwen3-Omni-Instruct (Xu et al., 2025b) | Qwen/Qwen3-Omni-30B-A3B-Instruct | model card | 30B (A3B) |
| Ming-Lite-Omni-1.5 (AI et al., 2025a) | inclusionAI/Ming-Lite-Omni-1.5 | model card | 20B (A3B) |
| Gemini-2.5-Flash (Comanici et al., 2025) | gemini-2.5-flash (June 17, 2025) | model card | - |

Table 6: Benchmark Candidates.

GLM-4-Voice. (Zeng et al., 2024) Co-developed by Zhipu.AI and Tsinghua University, it enables real-time voice conversations, and adjusts vocal nuances (such as emotion, intonation, and speech rate) according to user instructions. It also adopts a 12.5Hz single-codebook speech tokenizer and is pre-trained based on GLM-4-9B model, achieving SOTA performance across diverse tasks.

Salmonn. (Tang et al., 2024) Jointly developed by Tsinghua University and ByteDance, it integrates a pre-trained LLM (Vicuna) with dual auditory encoders (Whisper speech encoder and BEATs audio encoder). Salmonn uses a window-level Q-Former for cross-modal alignment and LoRA for

LLM adaptation. It achieves competitive performance on trained tasks and emergent abilities such as speech translation for untrained languages and audio-based storytelling.

MU-LLaMA. (Liu et al., 2023) Developed by Tencent ARC Lab and the National University of Singapore, MU-LLaMA is built on LLaMA, using a pretrained MERT model as the music encoder and trained on the specially constructed MusicQA dataset. It fuses music features into the LLaMA model via a Music Understanding Adapter, enabling both music-related question answering and music caption generation.

BAT. (Zheng et al., 2025b) Co-developed by the University of Texas at Austin and Shanghai Jiao Tong University, it integrates a novel spatial audio encoder called SPATIAL-AST (which excels in sound event detection, spatial localization, and distance estimation) with the LLaMA-2 7B, enabling it to perceive and reason about spatial sounds in 3D environments.

Step-Audio-2. (Wu et al., 2025) Developed by StepFun Audio Team, it’s tailored for industry-strength audio understanding and speech conversation. It adopts a unified architecture integrating a frozen latent audio encoder, an audio adaptor, an LLM decoder, and an audio detokenizer (Flow Matching + HiFi-GAN vocoder). It incorporates discrete audio token generation into language modeling, supports retrieval-augmented generation (RAG), and enables calling external tools (web search, audio search) to mitigate hallucination and switch timbres. It achieves state-of-the-art performance across multiple tasks.

H.2 Large Audio Reasoning Models

Mellow. (Deshmukh et al., 2025) Developed by CMU, it’s a small Audio Language Model tailored for audio-text reasoning tasks. It combines the HTSAT audio encoder and SmolLM2 small language model, trained on the ReasonAQA dataset. With only 167M parameters, it uses 50 times fewer parameters and 60 times less training audio than larger models while matching Qwen2 Audio’s performance on the MMAU benchmark and outperforming many larger models in deductive and comparative reasoning.

Audio-Reasoner. (Xie et al., 2025) Jointly developed by Nanyang Technological University, Skywork AI, Beijing Institute of Technology, and the National University of Singapore. It’s based on Qwen2-Audio-Instruct, fine-tuned on the 1.2-million-sample CoTA dataset via structured chain-of-thought (CoT) training to enhance deep audio reasoning capabilities. It achieves SOTA performance across key benchmarks.

GAMA. (Ghosh et al., 2024) Developed by UMD and Adobe, it integrates an LLM with multiple audio representations—including a custom Audio Q-Former and an Audio Spectrogram Transformer (AST) equipped with a multi-layer aggregator—fine-tuned on a large-scale audio-language dataset. It’s further instruction-tuned on

the synthetic CompA-R dataset to enhance complex reasoning, with high-level semantic evidence from audio event tags added via soft prompts.

R1-AQA. (Li et al., 2025a) This model is developed by Xiaomi Corporation, an enhanced version of Qwen2-Audio-7B-Instruct optimized via the Group Relative Policy Optimization (GRPO) for audio question answering tasks. It achieves SOTA performance on the MMAU Test-mini benchmark using only 38k post-training samples from the AVQA dataset, outperforming supervised fine-tuning methods even with its 8.2B parameters.

Audio Flamingo 2 Sound-CoT (Kong et al., 2025) Developed by NVIDIA, this work encompasses CoT-enhanced versions of Audio Flamingo 2 and Audio Flamingo 3. They share the same AF-CoT-Train dataset with 1.24M samples, which are constructed through four interactive pipelines between LLMs and ALMs to ensure audio-specific reasoning. The 7B-parameter Audio Flamingo 3 Sound-CoT sets a new state-of-the-art (SOTA) on the MMAU-Sound benchmark, delivers notable gains on AF-Reasoning-Eval (Classification) and MMAR-Sound, and maintains strong performance in discriminating closely related sound categories.

Audio Flamingo 3. (Goel et al., 2025) A collaboration between NVIDIA and UMD, it’s an LALM that advances reasoning and understanding across speech, sound, and music. It integrates AF-Whisper, adopts a five-stage curriculum-based training strategy, and supports key capabilities like multi-turn multi-audio chat, on-demand chain-of-thought reasoning, 10-minute long audio understanding, and voice-to-voice interaction. It leverages four novel datasets (AudioSkills-XL, LongAudio-XL, AF-Think, AF-Chat), achieves SOTA results on over 20 audio benchmarks.

Step-Audio-R1. (Tian et al., 2025) An open audio reasoning LLM from StepFun that explicitly targets the inverted scaling anomaly in audio models, where longer CoT previously degraded accuracy. Built on a frozen Qwen2 audio encoder and a Qwen2-style language backbone with an audio adaptor, Step-Audio-R1 is trained with Modality-Grounded Reasoning Distillation (MGRD), an iterative framework that filters and distills reasoning traces grounded in acoustic cues.

MiMo-Audio. (Xiaomi, 2025) A reasoning audio model developed by Xiaomi, its architec-

ture includes the MiMo-Audio-Tokenizer, a 1.2B-parameter Transformer operating at 25 Hz. MiMo-Audio is trained from scratch on a 10-million-hour corpus with joint optimization of semantic and reconstruction objectives, the tokenizer achieves superior reconstruction quality. It supports diverse tasks such as speech-to-speech generation, TTS, audio understanding, and spoken/text dialogue, demonstrating strong few-shot learning abilities.

H.3 Omni Language Models.

Qwen2.5-Omni. (Xu et al., 2025a) From Qwen2.5 series, developed by the Qwen Team, it’s an end-to-end multimodal model that perceives text, images, audio, and video. It adopts innovative techniques like TMRoPE for audio-video timestamp synchronization and the Thinker-Talker architecture to avoid text-speech interference, block-wise processing for multimodal encoders and a sliding-window DiT for low-latency audio streaming. It matches Qwen2.5-VL in image capabilities, outperforms Qwen2-Audio in audio tasks, achieves SOTA results on various benchmarks.

Baichuan-Omni-1.5. (Li et al., 2025c) Developed by Baichuan Inc., Baichuan-Omni is built on a high-quality dataset of about 500B samples, a custom 8-layer RVQ Baichuan-Audio-Tokenizer, and a multi-stage training strategy. It outperforms leading open-source models across text, image, video, and audio benchmarks, also achieves SOTA results on medical benchmarks.

VITA. (Fu et al., 2025) The omni version of VITA-audio, it’s developed by researchers from multiple institutions. This model starts with Mixtral 8×7B, expands its Chinese vocabulary, and through multimodal alignment and instruction tuning, enables processing of video, image, text, and audio. It features non-awakening and audio interrupt interactions.

Ola. (Liu et al., 2025c) Developed by Tsinghua University, Tencent Hunyuan Research, and S-Lab (NTU), Ola is an OLM built on the Qwen2.5-7B model. It integrates advanced encoders (OryxViT for vision, Whisper-v3 and BEATs for audio) and a progressive modality alignment strategy (starting with text-image, then adding video, and finally bridging vision-audio via cross-modal video data). It achieves competitive performance across image, video, and audio tasks, outperforming existing open omni-modal models.

HumanOmni. (Zhao et al., 2025) It’s the industry’s first human-centric vision-speech large language model, featuring three specialized branches (face-related, body-related, interaction-related) that adaptively fuse features via user instructions. It is trained on a dataset of over 2.4 million human-centric video clips and 14 million instructions, achieving state-of-the-art performance in tasks like emotion recognition, facial expression description, and action understanding.

OpenOmni. (Luo et al., 2025) It is an open-source omni language model developed by researchers from multiple institutions, aiming to address the scarcity of high-quality open omnimodal datasets and the challenge of real-time emotional speech synthesis. It adopts a two-stage framework for omnimodal alignment and speech generation, enabling vision-to-speech generalization and real-time emotional speech synthesis. It achieves competitive performance on multiple omnimodal benchmarks with a compact model size.

Qwen3-Omni. (Xu et al., 2025c) From the Qwen series, developed by the Qwen Team, it is a native end-to-end multilingual omnimodal foundation model and supports real-time streaming responses in both text and natural speech. It adopts a MoE-based Thinker-Talker architecture with AuT pretraining and a multi-codebook design to achieve strong cross-modal representations while minimizing latency. Qwen3-Omni achieves SOTA performance on a wide range of audio and video benchmarks.

Ming-Lite-Omni-1.5. (AI et al., 2025b) Ming-Omni adopts modality-specific encoders and an MoE-based core model, Ling, equipped with newly designed modality-aware routers to efficiently fuse multimodal inputs within a single framework, enabling diverse tasks without task-specific fine-tuning or architectural redesign. By integrating an advanced audio decoder for natural speech synthesis and Ming-Lite-Uni for high-quality image generation, Ming-Omni extends beyond perception to unified multimodal generation, and is the first open-source model to match GPT-4o in modality coverage.

H.4 Proprietary Models

GPT-4o-Audio-(2025-06-03). (OpenAI, 2024) Developed by OpenAI. It is a speech-centric iteration of the GPT-4o omni-modal family, designed to

handle native audio input and output directly without external ASR or TTS pipelines. This specific snapshot (2025-06-03) is optimized for low-latency, turn-based vocal interactions and agentic tasks, featuring enhanced instruction-following capabilities for complex audio analysis and generation.

Gemini-2.5-Flash (June 17, 2025). (Comanici et al., 2025) Developed by Google DeepMind. As a high-efficiency model within the Gemini 2.5 family, it integrates hybrid reasoning architectures with native multimodal understanding (text, audio, and video). It distinguishes itself through exceptional inference speed and cost-effectiveness while maintaining a large context window (up to 1M tokens), making it particularly suitable for real-time, long-context audio-visual reasoning tasks.

I Answer Extraction Details

Although our prompts explicitly instruct models to answer with a single letter (“A” or “B”), generative models often produce verbose responses or include conversational filler. To strictly parsing the outputs while maintaining fairness, we employ a cascading heuristic extraction strategy.

The extraction pipeline operates in the following order:

1. **Exact Match:** The output is stripped of whitespace and converted to lowercase. If it matches “a” or “b”, it is accepted.
2. **Normalization:** Common punctuation characters (e.g., brackets, periods, colons, hyphens) are replaced with whitespace. If the cleaned string is “a” or “b”, it is accepted.
3. **Pattern Matching:** If the direct methods fail, we iterate through the following set of regular expressions (case-insensitive) in order. The first pattern to yield a match determines the final prediction.

The specific regular expressions used are listed below:

1. `^s*([ab])s*[.]?s*$` # Matches “A”, “a.”, “b)”
2. `\boption\s*[:\~]?s*([ab])\b` # Matches “Option: A”
3. `\banswer\s*[:\~]?s*([ab])\b` # Matches “Answer: a”
4. `^s*\(?s*([ab])s*\)?s*$` # Matches “(A)”
5. `\b([ab])\b` # Fallback: distinct “A” or “B”

J Detailed Performance Breakdown

In the Section 4.4, we reported the aggregated accuracy across attributes to highlight the overall headroom for SOTA models. To facilitate a granular

analysis of model capabilities, we provide the full breakdown of accuracy scores separated by task type in this section.

Table 13 and Table 14 present the performance of all 36 evaluated systems on the **Comparison** and **Recognition** tasks across the 12 physical attributes respectively.

K Extended Analysis on Comparison vs. Recognition

In Section 4.4, we observed that current models do not exhibit the “comparison advantage” typical of human perception, where relative judgment (comparison) is generally easier than absolute estimation (recognition). Here, we provide supplementary data and visualizations to substantiate this finding.

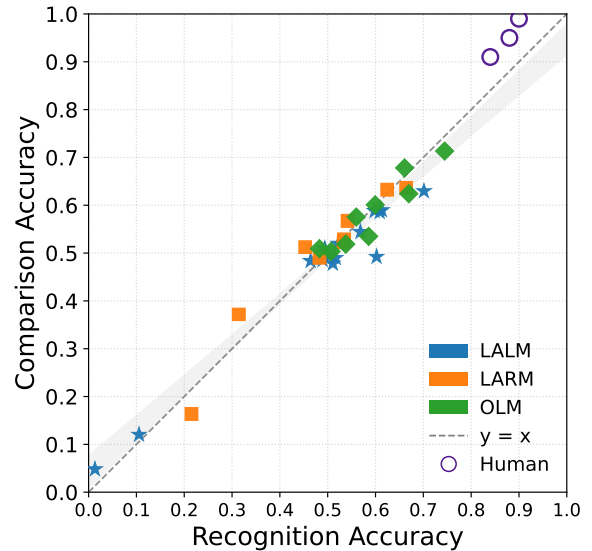


Figure 13: **Parity of model accuracy between comparison and recognition tasks.** Each point represents a model’s mean accuracy on comparison and recognition tasks across all attributes. The strong linear correlation ($r = 0.97$) indicates that current systems perform almost identically on the two task types, showing no systematic advantage for comparison-unlike humans, who typically benefit from contrastive cues between paired sounds.

K.1 Human Baselines

Table 7 details the performance of the three human participants involved in our study. As hypothesized and supported by psychophysical literature, human participants consistently achieve higher or equal accuracy in Comparison tasks compared to Recognition tasks across nearly all attributes. This advantage stems from the availability of contrastive cues

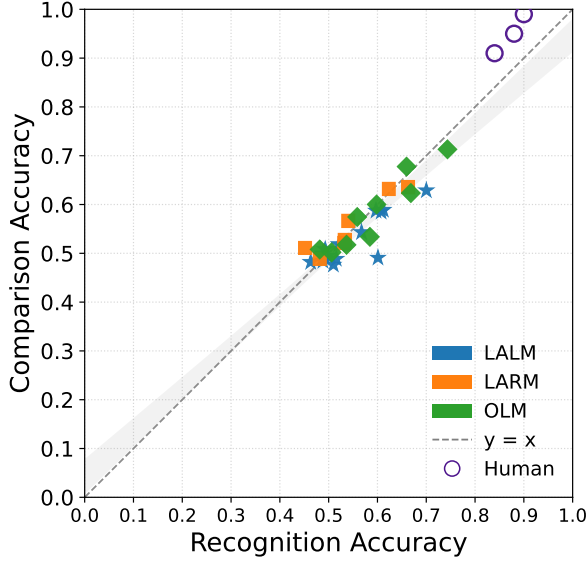


Figure 14: **Robustness of task parity.** This plot replicates the analysis in Figure 13 but excludes the four models with the highest abstention rates. The correlation remains unchanged ($r = 0.97$), confirming that the lack of comparison advantage is a fundamental characteristic of the models, not an artifact of poor instruction following.

between paired stimuli, which simplify decision-making by anchoring the judgment to an immediate reference.

K.2 Model Task Parity and Gaps

In contrast to humans, evaluated models typically show no such systematic advantage. We illustrate this through two complementary visualizations: global parity analysis and family-level gap analysis.

Global Parity. Figure 13 plots the mean accuracy of each model on recognition versus comparison tasks. We observe a tight linear correlation ($r = 0.97$), indicating that model performance is almost identical across task types. This suggests that current systems likely rely on similar internal mechanisms for both tasks, failing to exploit the pairwise contrastive information that benefits humans. To ensure this trend is not skewed by models with high failure rates, Figure 14 repeats this analysis excluding the four models with the highest abstention rates; the correlation remains robust ($r=0.97$), confirming the universality of this phenomenon.

Family-level Gaps. To further dissect these patterns, Figure 16 presents a dumbbell chart visualizing the accuracy gap between comparison and recognition tasks for each model family. Unlike the

| Tasks | Comparison Recognition | | |
|---------------|-----------------------------|-----------------------------|-----------------------------|
| | Participant 1 | Participant 2 | Participant 3 |
| Pitch | 1.00 0.80 ^{+25.0%} | 1.00 0.90 ^{+11.1%} | 1.00 0.90 ^{+11.1%} |
| Brightness | 1.00 1.00 ^{0.0%} | 0.80 0.80 ^{0.0%} | 1.00 1.00 ^{0.0%} |
| Loudness | 0.90 0.70 ^{+28.6%} | 1.00 0.80 ^{+25.0%} | 1.00 0.80 ^{+25.0%} |
| Velocity | 0.90 1.00 ^{-10.0%} | 0.60 0.70 ^{-14.3%} | 1.00 0.80 ^{+25.0%} |
| Duration | 1.00 0.70 ^{+42.9%} | 1.00 1.00 ^{0.0%} | 1.00 0.80 ^{+25.0%} |
| Tempo | 0.90 0.80 ^{+12.5%} | 1.00 0.70 ^{+42.9%} | 1.00 0.80 ^{+25.0%} |
| Direction | 0.80 0.90 ^{-11.1%} | 0.80 0.70 ^{+14.3%} | 0.90 0.90 ^{0.0%} |
| Distance | 0.90 0.70 ^{+28.6%} | 0.70 0.70 ^{0.0%} | 1.00 0.80 ^{+25.0%} |
| Reverberation | 1.00 1.00 ^{0.0%} | 1.00 1.00 ^{0.0%} | 1.00 1.00 ^{0.0%} |
| Texture | 1.00 1.00 ^{0.0%} | 1.00 1.00 ^{0.0%} | 1.00 1.00 ^{0.0%} |
| Timbre | 1.00 1.00 ^{0.0%} | 1.00 0.80 ^{+25.0%} | 1.00 1.00 ^{0.0%} |
| Counting | 1.00 1.00 ^{0.0%} | 1.00 1.00 ^{0.0%} | 1.00 1.00 ^{0.0%} |

Table 7: **Comparison vs. Recognition accuracy across acoustic attributes for three human Participants.** Each cell shows the mean accuracy for a given attribute pair (C | R), followed by the relative delta $\Delta\% = \frac{C-R}{R} \times 100$. Deltas are color-coded within a $\pm 5\%$ band, ($>5\%$) ($<-5\%$) ($|\leq 5\%$). Complementary with Figure 13, this table highlights that humans show consistent advantage on comparison tasks across attributes.

consistent positive gap observed in humans (Comparison > Recognition), model families exhibit inconsistent behaviors. While some OLMs show a slight preference for comparison, many LALM and LARM families exhibit negligible gaps or even inverted performance (where Recognition > Comparison), further highlighting the disconnect between biological hearing mechanisms and current model architectures.

L Case Study

In Section 4.4, we discussed the limitations of inference-time scaling for physical perception tasks. While Figure 4 in the main text demonstrates that reasoning often fails to rectify upstream perceptual errors, we also observe a more concerning phenomenon where explicit reasoning actively degrades performance by overriding correct perceptual intuitions.

Figure 15 provides a qualitative example of this reasoning-induced error on a reverberation comparison task. Notably, the model in its base mode correctly identifies the target audio (Option B), suggesting that the underlying acoustic representation implicitly includes the necessary physical cues. However, when the reasoning (think) mode is enabled, the model generates a verbose chain-of-thought. Although this reasoning chain is logically structured, it hallucinates incorrect acoustic descriptions that contradict the raw signal. This process effectively decouples the final prediction

from the acoustic reality, leading the model to reject the correct option and confidently select the wrong answer (Option A). This highlights a critical risk: without robust grounding, the CoT process could drift into hallucination, turning a correct perceptual detection into a reasoning failure.

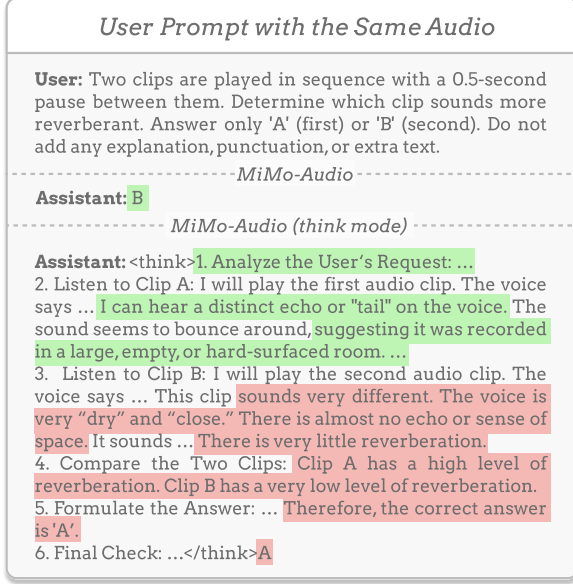


Figure 15: **A Case of reasoning-induced perceptual errors.** Both receive the same reverberation comparison prompt. The model in base mode correctly answers B, while in think mode produces a logically coherent reasoning but introduces erroneous explanation and answers A.

M Attribute Coverage Analysis

In addition to the per-attribute difficulty analysis presented in the Section 4.5, we evaluate the breadth of each model’s perceptual competence. We define “attribute coverage” as the number of attributes for which a model achieves a robust accuracy of ≥ 0.60 . This metric serves as a proxy for a model’s holistic understanding of the auditory world, distinguishing between models that excel in isolated tasks and those with generalized perceptual abilities.

Figure 6 visualizes this coverage distribution. The results highlight a distinct hierarchy in perceptual breadth: OLMs consistently demonstrate the widest coverage, led by Qwen3-Omni (9 attributes). LARMs occupy the middle ground, while LALMs generally exhibit the narrowest scope. Notably, no model in the LALM category meets the reference threshold of 7.2 attributes (representing 60% coverage of the 12 total attributes).

N Probe Experiment Details

N.1 Probing Model Details

We detail the pipeline components and weight initialization sources for each E2E model in Table 8. Encoder Status indicates whether the audio encoder parameters were updated during the model’s specific instruction-tuning or alignment stage.

N.2 Probe Architecture

To verify the intrinsic perceptual quality of the representations, we employ a lightweight probing architecture. The probe consists of the pre-trained audio encoder or fine-tuned audio encoder extracted from each respective LALM, followed by a shallow classification head. Formally, given an input audio x , the frozen encoder produces a sequence of hidden states $\mathbf{H} \in \mathbb{R}^{T \times D_{enc}}$, where T is the sequence length and D_{enc} is the hidden dimension.

The probing head consists of:

1. A trainable linear projection layer transforming features from D_{enc} to a dimension $D_{proj} = 256$.
2. A mean pooling layer to aggregate the sequence into a single vector.
3. A linear classification layer mapping the pooled vector to the label space with **2 output classes** (binary classification).

During training, the parameters of the audio encoder are strictly **frozen**, and only the projection and classification layers are updated.

N.3 Training Hyperparameters

We train all probes using the AdamW optimizer with a cosine learning rate schedule. To prevent overfitting, we employ an early stopping mechanism based on validation accuracy. The unified training configurations across all experiments are detailed in Table 9.

O Detailed Probe Results

In Section 4.6, we summarized the overall bottleneck analysis using aggregated accuracy. To provide a granular view of where this information loss occurs, Table 10 details the accuracy for both Recognition and Comparison tasks across all 12 attributes.

The table compares three model states for each of the eight representative systems:

| E2E Model | Audio Encoder Arch. | Projector | LLM Backbone | Encoder Init. | LLM Init. | Encoder Status |
|-------------------------|--|---|-------------------------------------|--|-----------------|----------------|
| SALMONN | BEATs (Chen et al., 2022) & Whisper-Large-v2 | Window-level Q-Former (Li et al., 2023) | Vicuna-7B-v1.5 (Zheng et al., 2023) | Fine-tuned_BEATs_iter3+(AS2M)(cpt2) & Whisper-Large-v2 | Vicuna-7B-v1.5 | Frozen |
| Step-Audio-2-mini | Qwen2-Audio Encoder (Whisper-Large-v3) | Downsampling Adaptor | Qwen2.5-7B (Team, 2024) | Qwen2-Audio | Qwen2.5-7B | Frozen |
| VITA-Audio-Plus-Vanilla | SenseVoiceSmall (An et al., 2024) | MLP | Qwen2.5-7B | SenseVoiceSmall | Qwen2.5-7B | Frozen |
| MiDashengLM | Dasheng-0.6B (Dinkel et al., 2024) | MLP | Qwen2.5-Omni-7B | Dasheng-0.6B | Qwen2.5-Omni-7B | Unfrozen |
| Qwen2-Audio | Whisper-Large-v3 (Radford et al., 2022) | - | Qwen2.5-7B (Bai et al., 2023) | Whisper-Large-v3 | Qwen2.5-7B | Unfrozen |
| Kimi-Audio | Whisper-Large-v3 | Downsampling Adaptor | Qwen2.5-7B | Whisper-Large-v3 | Qwen2.5-7B | Unfrozen |
| Qwen2.5-Omni | Whisper-Large-v3 | - | Qwen2.5-7B | Whisper-Large-v3 | Qwen2.5-7B | Unfrozen |
| Qwen3-Omni | AuT (Xu et al., 2025c) | - | Qwen3 | AuT | Qwen3 | Unfrozen |

Table 8: Architecture and Training Configurations of Evaluated Models.

| Hyperparameter | Value |
|-------------------------------------|--------------------|
| Optimizer | AdamW |
| Learning Rate | 3×10^{-3} |
| LR Scheduler | Cosine Decay |
| Warmup Ratio | 0.1 |
| Weight Decay | 0.01 |
| Adam β_1, β_2 | 0.9, 0.98 |
| Adam ϵ | 1×10^{-8} |
| Batch Size (per device) | 2 |
| Training Epochs | 20 |
| Early Stopping Patience | 3 epochs |
| Early Stopping Threshold | 0.001 |
| Projection Dimension (D_{proj}) | 256 |
| Number of Classes | 2 |
| Seed | 42 |

Table 9: **Hyperparameters for Probe Training.** These settings are consistent across all model families and tasks.

1. **Pretrained Encoder (Probe):** The raw capability of the frozen encoder before instruction tuning.
2. **Extracted Encoder (Probe):** The capability of the encoder after end-to-end training (only for models that unfreeze the encoder).
3. **E2E Model (Zero-shot):** The final performance of the full multimodal system.

The substantial performance drop (Δ) from Probe to E2E across high-performing perceptual attributes (e.g., Pitch, Brightness) serves as direct evidence of the alignment/decoding bottleneck.

O.1 The Relational Bottleneck in Encoders

While the previous analysis highlighted the degradation from encoder to E2E model, we also investigate whether the encoders themselves possess the structural capability to handle relational tasks. Table 11 breaks down the linear probe performance results across the two tasks, Recognition (absolute judgment) and Comparison (relative judgment).

Contrary to human perception, where comparison is cognitively less demanding, our probing results reveal that encoders do not exhibit a “comparison advantage.” This suggests that standard audio

encoders, primarily trained on global or frame-level classification objectives, lack the inherent mechanisms to explicitly compare distinct temporal segments within a single audio stream, limiting their ability to leverage the contrastive structure of the input.

O.2 Attribute-wise Patterns in Encoders

Beyond identifying bottlenecks, our probing setup allows us to analyze the plasticity of audio encoders during the E2E training process. Table 12 compares the probing accuracy of *Pretrained Encoders* versus *Extracted Encoders*. We observe that E2E training with an unfrozen audio encoder can yield small gains.

To deeper inspection in Table 11, we find a non-uniform adaptation capability across different physical attributes. Comparing pre-trained encoders and extracted encoders, we can tell that unfreezing the encoder helps the Tempo and Distance attributes more noticeably, yet brings limited gains for Timbre, Texture, and Counting.

P Implications

Concretely, our findings suggest two distinct pathways for future development: (i) **For perceptual recognition.** Prioritize *alignment strategies*. Since encoders already encode rich signal cues as shown by our probes, future work should focus on preventing representational degradation during LLM integration, potentially via adaptable projectors or the partial unfreezing strategies observed in Omni-models. (ii) **For relational reasoning.** Innovate in *encoder pretraining*. Current encoders lack the mechanisms to explicitly compare inputs. Addressing this requires incorporating architectural inductive biases, such as native cross-segment attention or contrastive objectives, to ensure relational information is captured before reaching the LLM.

| Accuracy Across Model Groups | | | | | | | | | | | | | | | | | | | | | |
|------------------------------|-------|------------------|-------------------|---------------------------|-------------------|------------------|-------------------------|---------------------|-----------------------|-------------------|-------------------|------------------------------|----------------------|-------------------|-----------------------------|---------------------|-------------------|-------------------------------|-----------------------|-----------------------------|---------------------|
| (a) Comparison Accuracy | | | | | | | | | | | | | | | | | | | | | |
| Attributes | BEATs | Whisper-Large-v2 | SAL-MONN | Step-Audio 2 mini Encoder | Step-Audio 2 mini | SenseVoice Small | VITA-Audio-Plus-Vanilla | Da-sheng LM Encoder | MidDasheng-LM Encoder | MidDashengLM | Whisper-Large-v3 | Qwen2-Audio-Instruct Encoder | Qwen2-Audio-Instruct | Whisper-Large-v3 | Kimi-Audio-Instruct Encoder | Kimi-Audio-Instruct | Whisper-Large-v3 | Qwen2.5-Omni-Instruct Encoder | Qwen2.5-Omni-Instruct | Qwen3-Omni-Instruct Encoder | Qwen3-Omni-Instruct |
| Pitch | 0.66 | 0.76 | 0.50 ⁺ | 0.72 | 0.47 ⁺ | 0.66 | 0.52 ⁺ | 0.56 | 0.82 | 0.56 ⁺ | 0.74 | 0.60 | 0.43 ⁺ | 0.74 | 0.74 | 0.79 | 0.74 | 0.66 | 0.60 ⁺ | 0.58 | 0.83 |
| Brightness | 0.70 | 0.92 | 0.50 ⁺ | 0.70 | 0.55 ⁺ | 0.58 | 0.51 ⁺ | 0.62 | 0.76 | 0.64 ⁺ | 0.84 | 0.60 | 0.49 ⁺ | 0.84 | 0.74 | 0.81 ⁺ | 0.84 | 0.60 | 0.89 | 0.52 | 0.92 |
| Loudness | 0.58 | 0.58 | 0.53 ⁺ | 0.66 | 0.57 ⁺ | 0.62 | 0.50 ⁺ | 0.62 | 0.68 | 0.61 ⁺ | 0.58 | 0.66 | 0.49 ⁺ | 0.58 | 0.70 | 0.66 ⁺ | 0.58 | 0.62 | 0.73 | 0.60 | 0.80 |
| Velocity | 0.54 | 0.56 | 0.50 ⁺ | 0.60 | 0.48 ⁺ | 0.54 | 0.45 ⁺ | 0.56 | 0.52 | 0.62 | 0.56 | 0.56 | 0.56 | 0.56 | 0.58 | 0.53 ⁺ | 0.56 | 0.56 | 0.54 ⁺ | 0.50 | 0.64 |
| Duration | 0.52 | 0.56 | 0.49 ⁺ | 0.58 | 0.47 ⁺ | 0.50 | 0.51 | 0.58 | 0.56 | 0.58 | 0.50 | 0.56 | 0.52 ⁺ | 0.50 | 0.62 | 0.64 | 0.50 | 0.54 | 0.54 | 0.50 | 0.73 |
| Tempo | 0.48 | 0.58 | 0.52 ⁺ | 0.60 | 0.52 ⁺ | 0.54 | 0.49 ⁺ | 0.56 | 0.56 | 0.55 ⁺ | 0.62 | 0.60 | 0.55 ⁺ | 0.62 | 0.54 | 0.58 ⁺ | 0.62 | 0.52 | 0.50 ⁺ | 0.52 | 0.66 |
| Direction | 0.52 | 0.68 | 0.51 ⁺ | 0.54 | 0.51 ⁺ | 0.48 | 0.45 ⁺ | 0.48 | 0.46 | 0.48 | 0.54 | 0.52 | 0.52 ⁺ | 0.54 | 0.54 | 0.43 ⁺ | 0.54 | 0.52 | 0.50 | 0.44 | 0.52 |
| Distance | 0.58 | 0.62 | 0.47 ⁺ | 0.58 | 0.50 ⁺ | 0.56 | 0.54 ⁺ | 0.54 | 0.60 | 0.54 ⁺ | 0.56 | 0.52 | 0.54 ⁺ | 0.56 | 0.64 | 0.39 ⁺ | 0.56 | 0.52 | 0.52 ⁺ | 0.52 | 0.60 |
| Reverberation | 0.54 | 0.80 | 0.50 ⁺ | 0.68 | 0.48 ⁺ | 0.50 | 0.54 | 0.64 | 0.56 | 0.54 ⁺ | 0.84 | 0.56 | 0.56 ⁺ | 0.84 | 0.54 | 0.58 ⁺ | 0.84 | 0.52 | 0.42 ⁺ | 0.52 | 0.52 |
| Timbre | 0.48 | 0.54 | 0.50 ⁺ | 0.58 | 0.91 | 0.54 | 0.50 ⁺ | 0.60 | 0.62 | 0.74 | 0.54 | 0.66 | 0.50 ⁺ | 0.54 | 0.52 | 0.95 | 0.54 | 0.62 | 1.00 | 0.50 | 0.99 |
| Texture | 0.54 | 0.50 | 0.50 ⁺ | 0.50 | 0.50 | 0.58 | 0.50 ⁺ | 0.62 | 0.56 | 0.51 ⁺ | 0.50 | 0.58 | 0.50 ⁺ | 0.50 | 0.64 | 0.49 ⁺ | 0.50 | 0.52 | 0.69 | 0.54 | 0.70 |
| Counting | 0.56 | 0.50 | 0.50 ⁺ | 0.76 | 0.50 ⁺ | 0.70 | 0.48 ⁺ | 0.82 | 0.74 | 0.60 ⁺ | 0.50 | 0.56 | 0.59 | 0.50 | 0.84 | 0.79 ⁺ | 0.50 | 0.52 | 0.55 | 0.62 | 0.65 |
| Per-task Average | 0.56 | 0.63 | 0.50 ⁺ | 0.63 | 0.54 ⁺ | 0.57 | 0.50 ⁺ | 0.60 | 0.62 | 0.58 ⁺ | 0.61 | 0.58 | 0.52 ⁺ | 0.61 | 0.64 | 0.64 | 0.61 | 0.56 | 0.62 | 0.53 | 0.71 |
| Accuracy Across Model Groups | | | | | | | | | | | | | | | | | | | | | |
| (b) Recognition Accuracy | | | | | | | | | | | | | | | | | | | | | |
| Attributes | BEATs | Whisper-Large-v2 | SAL-MONN | Step-Audio 2 mini Encoder | Step-Audio 2 mini | SenseVoice Small | VITA-Audio-Plus-Vanilla | Da-sheng LM Encoder | MidDasheng-LM Encoder | MidDashengLM | Whisper-Large-v3 | Qwen2-Audio-Instruct Encoder | Qwen2-Audio-Instruct | Whisper-Large-v3 | Kimi-Audio-Instruct Encoder | Kimi-Audio-Instruct | Whisper-Large-v3 | Qwen2.5-Omni-Instruct Encoder | Qwen2.5-Omni-Instruct | Qwen3-Omni-Instruct Encoder | Qwen3-Omni-Instruct |
| Pitch | 0.92 | 0.92 | 0.50 ⁺ | 0.92 | 0.57 ⁺ | 0.76 | 0.50 ⁺ | 0.88 | 0.88 | 0.65 ⁺ | 0.86 | 0.96 | 0.52 ⁺ | 0.86 | 0.90 | 0.87 ⁺ | 0.86 | 0.92 | 0.82 ⁺ | 0.94 | 0.90 ⁺ |
| Brightness | 0.96 | 0.92 | 0.50 ⁺ | 1.00 | 0.62 ⁺ | 0.84 | 0.50 ⁺ | 0.88 | 0.98 | 0.75 ⁺ | 0.84 | 0.98 | 0.49 ⁺ | 0.84 | 0.94 | 0.81 ⁺ | 0.84 | 0.86 | 0.84 ⁺ | 0.88 | 0.83 ⁺ |
| Loudness | 0.76 | 0.70 | 0.50 ⁺ | 0.76 | 0.41 ⁺ | 0.62 | 0.50 ⁺ | 0.70 | 0.74 | 0.72 ⁺ | 0.70 | 0.82 | 0.47 ⁺ | 0.70 | 0.80 | 0.70 ⁺ | 0.70 | 0.80 | 0.76 ⁺ | 0.88 | 0.70 ⁺ |
| Velocity | 0.60 | 0.62 | 0.50 ⁺ | 0.64 | 0.53 ⁺ | 0.60 | 0.51 ⁺ | 0.58 | 0.66 | 0.46 ⁺ | 0.68 | 0.66 | 0.54 ⁺ | 0.68 | 0.58 | 0.50 ⁺ | 0.68 | 0.62 | 0.52 ⁺ | 0.70 | 0.49 ⁺ |
| Duration | 0.90 | 0.78 | 0.50 ⁺ | 0.94 | 0.54 ⁺ | 0.90 | 0.50 ⁺ | 0.94 | 0.94 | 0.68 ⁺ | 0.78 | 1.00 | 0.50 ⁺ | 0.78 | 1.00 | 0.75 ⁺ | 0.78 | 0.98 | 0.60 ⁺ | 0.94 | 0.80 ⁺ |
| Tempo | 0.50 | 0.50 | 0.50 ⁺ | 0.68 | 0.47 ⁺ | 0.66 | 0.51 ⁺ | 0.86 | 0.78 | 0.55 ⁺ | 0.52 ⁺ | 0.60 | 0.51 ⁺ | 0.52 ⁺ | 0.64 | 0.55 ⁺ | 0.52 ⁺ | 0.62 | 0.53 ⁺ | 0.58 | 0.64 |
| Direction | 0.72 | 0.68 | 0.50 ⁺ | 0.76 | 0.43 ⁺ | 0.70 | 0.47 ⁺ | 0.74 | 0.74 | 0.44 ⁺ | 0.68 | 0.76 | 0.48 ⁺ | 0.68 | 0.76 | 0.44 ⁺ | 0.68 | 0.74 | 0.49 ⁺ | 0.74 | 0.49 ⁺ |
| Distance | 0.56 | 0.56 | 0.44 ⁺ | 0.60 | 0.46 ⁺ | 0.64 | 0.53 ⁺ | 0.58 | 0.58 | 0.48 ⁺ | 0.56 | 0.62 | 0.52 ⁺ | 0.56 | 0.60 | 0.48 ⁺ | 0.56 | 0.66 | 0.52 ⁺ | 0.60 | 0.42 ⁺ |
| Reverberation | 1.00 | 1.00 | 0.50 ⁺ | 1.00 | 0.52 ⁺ | 0.96 | 0.50 ⁺ | 1.00 | 1.00 | 0.50 ⁺ | 1.00 | 0.98 | 0.50 ⁺ | 1.00 | 1.00 | 0.95 ⁺ | 1.00 | 0.98 | 0.69 ⁺ | 1.00 | 0.83 ⁺ |
| Timbre | 0.50 | 0.62 | 0.50 ⁺ | 0.52 | 0.98 | 0.56 | 0.54 ⁺ | 0.54 | 0.48 | 0.96 | 0.58 | 0.52 | 0.67 ⁺ | 0.58 | 0.56 | 0.99 | 0.58 | 0.54 | 0.97 | 0.58 | 0.98 |
| Texture | 0.52 | 0.50 | 0.50 ⁺ | 0.56 | 0.65 | 0.56 | 0.48 ⁺ | 0.52 | 0.66 | 0.69 | 0.50 | 0.56 | 0.60 | 0.50 | 0.52 | 0.68 | 0.50 | 0.54 | 0.80 | 0.56 | 0.75 |
| Counting | 0.54 | 0.58 | 0.50 ⁺ | 0.52 | 0.52 ⁺ | 0.52 | 0.40 ⁺ | 0.62 | 0.58 | 0.44 ⁺ | 0.58 | 0.52 | 0.47 ⁺ | 0.58 | 0.54 | 0.69 | 0.58 | 0.56 | 0.60 | 0.58 | 0.99 |
| Per-task Average | 0.71 | 0.70 | 0.50 ⁺ | 0.74 | 0.56 ⁺ | 0.69 | 0.50 ⁺ | 0.74 | 0.75 | 0.61 ⁺ | 0.69 | 0.75 | 0.52 ⁺ | 0.69 | 0.74 | 0.70 ⁺ | 0.69 | 0.74 | 0.68 ⁺ | 0.75 | 0.74 ⁺ |

Table 10: Detailed Accuracy Breakdown: Encoders (Probe) vs. E2E Models. We report accuracy across 12 attributes separated by Recognition and Comparison tasks. *Pretrained* denotes the off-the-shelf encoder; *Extracted* refers to the encoder state after E2E training (if unfrozen). The results confirm a widespread alignment bottleneck that linear probes on encoders consistently achieve robust accuracy (≥ 0.60) on signal-level attributes, whereas corresponding E2E models often degrade to near-random performance (≈ 0.50). Notable exceptions are the Qwen-Omni series, where E2E performance matches or exceeds the probe, indicating a successful preservation of perceptual information during their training.

| Tasks | | Comparison Recognition | | | | | | | | | | |
|----------------------|---------------|--------------------------|------------------|------------------|------------------|-----------|---------------------|---------------------------|------------------------------|-----------------------------|-------------------------------|-----------------------------|
| Metric | Accuracy | BEATs | Whisper-Large-v2 | Whisper-Large-v3 | SenseVoice Small | Dasheng | MiDashengLM Encoder | Step-Audio 2-mini Encoder | Qwen2-Audio-Instruct Encoder | Kimi-Audio-Instruct Encoder | Qwen2.5-Omni-Instruct Encoder | Qwen3-Omni-Instruct Encoder |
| Spectral& Amplitude | Pitch | 0.66 0.92 | 0.76 0.92 | 0.74 0.86 | 0.66 0.76 | 0.56 0.88 | 0.82 0.88 | 0.72 0.92 | 0.60 0.96 | 0.74 0.90 | 0.66 0.92 | 0.58 0.94 |
| | Brightness | 0.70 0.96 | 0.92 0.92 | 0.84 0.84 | 0.58 0.84 | 0.62 0.88 | 0.76 0.98 | 0.70 1.00 | 0.60 0.98 | 0.74 0.94 | 0.60 0.86 | 0.52 0.88 |
| | Loudness | 0.58 0.76 | 0.58 0.70 | 0.58 0.70 | 0.62 0.62 | 0.62 0.70 | 0.68 0.74 | 0.66 0.76 | 0.66 0.82 | 0.70 0.80 | 0.62 0.80 | 0.60 0.88 |
| | Velocity | 0.54 0.60 | 0.56 0.62 | 0.56 0.68 | 0.54 0.60 | 0.56 0.58 | 0.52 0.66 | 0.60 0.64 | 0.56 0.66 | 0.58 0.58 | 0.56 0.62 | 0.50 0.70 |
| Temporal | Duration | 0.52 0.90 | 0.56 0.78 | 0.50 0.78 | 0.50 0.90 | 0.58 0.94 | 0.56 0.94 | 0.58 0.94 | 0.56 1.00 | 0.62 1.00 | 0.54 0.98 | 0.50 0.94 |
| | Tempo | 0.48 0.50 | 0.58 0.50 | 0.62 0.52 | 0.54 0.66 | 0.56 0.86 | 0.56 0.78 | 0.60 0.68 | 0.60 0.60 | 0.54 0.64 | 0.52 0.62 | 0.52 0.58 |
| Spatial& Environment | Direction | 0.52 0.72 | 0.68 0.68 | 0.54 0.68 | 0.48 0.70 | 0.48 0.74 | 0.46 0.74 | 0.54 0.76 | 0.52 0.76 | 0.54 0.76 | 0.52 0.74 | 0.44 0.74 |
| | Distance | 0.58 0.56 | 0.62 0.56 | 0.56 0.56 | 0.56 0.64 | 0.54 0.58 | 0.60 0.58 | 0.58 0.60 | 0.52 0.62 | 0.64 0.60 | 0.52 0.66 | 0.52 0.60 |
| | Reverberation | 0.54 1.00 | 0.80 1.00 | 0.84 1.00 | 0.50 0.96 | 0.64 1.00 | 0.56 1.00 | 0.68 1.00 | 0.56 0.98 | 0.54 1.00 | 0.52 0.98 | 0.52 1.00 |
| Timbre | Timbre | 0.48 0.50 | 0.54 0.62 | 0.54 0.58 | 0.54 0.56 | 0.60 0.54 | 0.62 0.48 | 0.58 0.52 | 0.66 0.52 | 0.52 0.56 | 0.62 0.54 | 0.50 0.58 |
| | Texture | 0.54 0.52 | 0.50 0.50 | 0.50 0.50 | 0.58 0.56 | 0.62 0.52 | 0.56 0.66 | 0.50 0.56 | 0.58 0.56 | 0.64 0.52 | 0.52 0.54 | 0.54 0.56 |
| Scene Level | Counting | 0.56 0.54 | 0.50 0.58 | 0.50 0.58 | 0.70 0.52 | 0.82 0.62 | 0.74 0.58 | 0.76 0.52 | 0.56 0.52 | 0.84 0.54 | 0.52 0.56 | 0.62 0.58 |
| Per-task Average | | 0.56 0.71 | 0.63 0.70 | 0.61 0.69 | 0.57 0.69 | 0.60 0.74 | 0.62 0.75 | 0.63 0.74 | 0.58 0.75 | 0.64 0.74 | 0.56 0.74 | 0.53 0.75 |
| Overall Average | | 0.63 | 0.67 | 0.65 | 0.63 | 0.67 | 0.69 | 0.68 | 0.67 | 0.69 | 0.65 | 0.64 |

Table 11: Probe Accuracy across Task Dimensions: Evidence of a Relational Bottleneck. This table reports the mean accuracy of linear probes trained on frozen pretrained encoders and their unfrozen encoders, aggregated by task type (Recognition vs. Comparison) across all 12 attributes. Crucially, comparison performance consistently remains at or below recognition performance, mirroring the behavior of full E2E models. This indicates a *shared relational bottleneck* that without architectural mechanisms to explicitly model contrastive relationships between segments, even raw perceptual representations fail to leverage the comparative structure of the task, treating comparison as merely a more complex form of recognition.

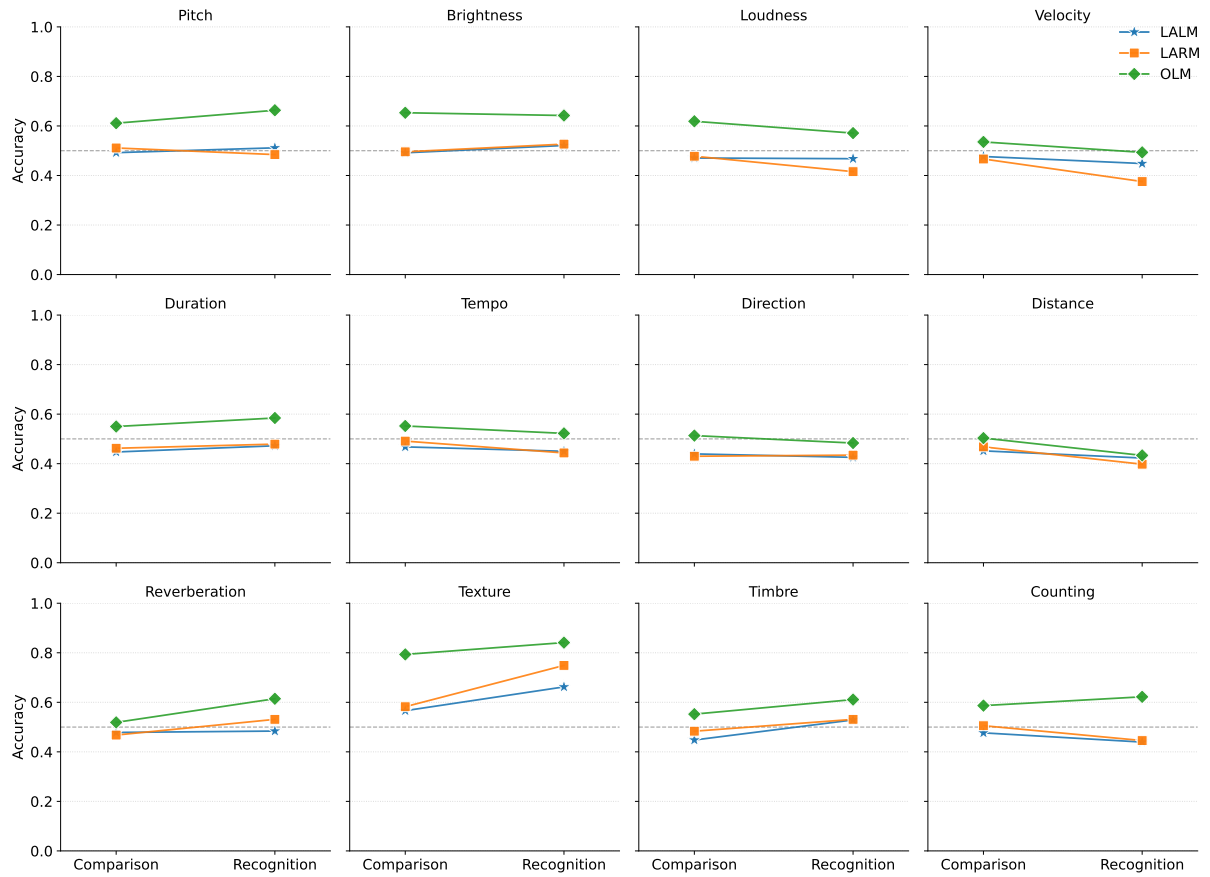


Figure 16: **Family-level performance gaps between comparison and recognition tasks.** This dumbbell chart illustrates the mean accuracy difference between the two task types for each model family. The two endpoints of each dumbbell represent the accuracy for Recognition and Comparison, respectively. Unlike the consistent “comparison advantage” observed in humans, model families show varied behaviors where gaps are generally narrow, and in several cases (especially within LALMs), comparison accuracy unexpectedly falls below recognition accuracy.

| Model Groups | Comparison Accuracy | Recognition Accuracy | Overall Accuracy |
|-------------------------------|---------------------|----------------------|-------------------|
| BEATs | 0.56 | 0.71 | 0.63 |
| Whisper-Large-v2 | 0.63 | 0.70 | 0.67 |
| Salmonn | 0.50 [↓] | 0.50 ⁺ | 0.50 [↓] |
| Whisper-Large-v3 | 0.61 | 0.69 | 0.65 |
| Step-Audio 2 mini Encoder | 0.63 | 0.74 | 0.68 |
| Step-Audio 2 mini | 0.54 [↓] | 0.56 ⁺ | 0.55 [↓] |
| SenseVoiceSmall | 0.57 | 0.69 | 0.63 |
| VITA-Audio-Plus-Vanilla | 0.50 [↓] | 0.50 ⁺ | 0.50 [↓] |
| Dasheng | 0.60 | 0.74 | 0.67 |
| MiDashengLM Encoder | 0.62 | 0.75 | 0.69 |
| MiDashengLM | 0.58 [↓] | 0.61 ⁺ | 0.60 [↓] |
| Whisper-Large-v3 | 0.61 | 0.69 | 0.65 |
| Qwen2-Audio-Instruct Encoder | 0.58 | 0.75 | 0.67 |
| Qwen2-Audio-Instruct | 0.52 [↓] | 0.52 ⁺ | 0.52 [↓] |
| Whisper-Large-v3 | 0.61 | 0.69 | 0.65 |
| Kimi-Audio-Instruct Encoder | 0.64 | 0.74 | 0.69 |
| Kimi-Audio-Instruct | 0.64 | 0.70 ⁺ | 0.67 [↓] |
| Whisper-Large-v3 | 0.61 | 0.69 | 0.65 |
| Qwen2.5-Omni-Instruct Encoder | 0.56 | 0.74 | 0.65 |
| Qwen2.5-Omni-Instruct | 0.62 | 0.68 ⁺ | 0.65 |
| Qwen3-Omni-Instruct Encoder | 0.53 | 0.75 | 0.64 |
| Qwen3-Omni-Instruct | 0.71 | 0.74 ⁺ | 0.72 |

Table 12: This table contrasts the performance of off-the-shelf Pretrained Encoders against Extracted Encoders that were updated during training and E2E models. We can find that E2E training with an unfrozen audio encoder can yield small gains

| Models | Spectral & Amplitude | | | | Temporal | | Spatial& Environment | | | Timbre | | Scene Level | Avg. |
|--------------------------------------|----------------------|------------|----------|----------|----------|-------|----------------------|----------|---------------|---------|--------|-------------|------|
| | Pitch | Brightness | Loudness | Velocity | Duration | Tempo | Direction | Distance | Reverberation | Texture | Timbre | Counting | |
| Random Guess | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 |
| Human | 1.00 | 0.93 | 0.97 | 0.83 | 1.00 | 0.97 | 0.83 | 0.87 | 1.00 | 1.00 | 1.00 | 1.00 | 0.95 |
| Large Audio Language Models (LALMs) | | | | | | | | | | | | | |
| BAT | 0.05 | 0.00 | 0.04 | 0.08 | 0.03 | 0.00 | 0.00 | 0.00 | 0.00 | 0.27 | 0.09 | 0.00 | 0.05 |
| MU-LLaMA | 0.34 | 0.08 | 0.00 | 0.32 | 0.04 | 0.12 | 0.10 | 0.16 | 0.24 | 0.00 | 0.00 | 0.06 | 0.12 |
| LLaMA-Omni | 0.38 | 0.50 | 0.51 | 0.50 | 0.46 | 0.51 | 0.50 | 0.50 | 0.49 | 0.50 | 0.50 | 0.44 | 0.48 |
| Audio Flamingo 2 | 0.50 | 0.47 | 0.53 | 0.50 | 0.49 | 0.52 | 0.45 | 0.39 | 0.49 | 0.48 | 0.50 | 0.50 | 0.49 |
| GLM-4-Voice | 0.50 | 0.51 | 0.47 | 0.49 | 0.47 | 0.54 | 0.44 | 0.50 | 0.49 | 0.50 | 0.50 | 0.46 | 0.49 |
| VITA-Audio-Plus-Vanilla | 0.52 | 0.51 | 0.50 | 0.45 | 0.51 | 0.49 | 0.45 | 0.54 | 0.54 | 0.50 | 0.50 | 0.48 | 0.50 |
| Audio Flamingo 3 | 0.50 | 0.50 | 0.50 | 0.50 | 0.44 | 0.52 | 0.51 | 0.52 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 |
| Voxtral-Mini | 0.53 | 0.45 | 0.53 | 0.51 | 0.49 | 0.52 | 0.51 | 0.47 | 0.50 | 0.52 | 0.50 | 0.52 | 0.50 |
| Baichuan-Audio-Instruct | 0.48 | 0.52 | 0.49 | 0.48 | 0.53 | 0.51 | 0.45 | 0.47 | 0.50 | 0.54 | 0.50 | 0.49 | 0.50 |
| LLaMA-Omni2 | 0.49 | 0.47 | 0.42 | 0.47 | 0.48 | 0.49 | 0.56 | 0.56 | 0.55 | 0.54 | 0.50 | 0.50 | 0.50 |
| SALMONN | 0.50 | 0.50 | 0.53 | 0.50 | 0.49 | 0.52 | 0.51 | 0.47 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 |
| R1-AQA | 0.50 | 0.49 | 0.50 | 0.50 | 0.46 | 0.51 | 0.50 | 0.48 | 0.52 | 0.51 | 0.48 | 0.49 | 0.50 |
| Qwen2-Audio-Instruct | 0.43 | 0.49 | 0.49 | 0.56 | 0.52 | 0.55 | 0.52 | 0.54 | 0.56 | 0.50 | 0.50 | 0.59 | 0.52 |
| Step-Audio 2 mini | 0.47 | 0.55 | 0.57 | 0.48 | 0.47 | 0.52 | 0.51 | 0.50 | 0.48 | 0.91 | 0.50 | 0.50 | 0.54 |
| MiDashengLM | 0.56 | 0.64 | 0.61 | 0.62 | 0.58 | 0.55 | 0.48 | 0.54 | 0.54 | 0.74 | 0.51 | 0.60 | 0.58 |
| MiMo-Audio-Instruct | 0.61 | 0.68 | 0.57 | 0.54 | 0.51 | 0.48 | 0.50 | 0.54 | 0.63 | 0.85 | 0.56 | 0.50 | 0.58 |
| Kimi-Audio-Instruct | 0.79 | 0.81 | 0.66 | 0.53 | 0.64 | 0.58 | 0.43 | 0.39 | 0.58 | 0.95 | 0.49 | 0.79 | 0.64 |
| GPT-4o-Audio | 0.71 | 0.68 | 0.55 | 0.55 | 0.44 | 0.49 | 0.49 | 0.56 | 0.50 | 0.88 | 0.43 | 0.66 | 0.58 |
| Large Audio Reasoning Models (LARMs) | | | | | | | | | | | | | |
| Mellow | 0.07 | 0.08 | 0.14 | 0.20 | 0.18 | 0.28 | 0.16 | 0.20 | 0.28 | 0.16 | 0.03 | 0.15 | 0.16 |
| GAMA | 0.38 | 0.42 | 0.38 | 0.45 | 0.24 | 0.49 | 0.29 | 0.36 | 0.26 | 0.48 | 0.41 | 0.33 | 0.37 |
| Audio Flamingo 2 Sound-CoT | 0.52 | 0.50 | 0.55 | 0.49 | 0.50 | 0.52 | 0.51 | 0.47 | 0.49 | 0.50 | 0.51 | 0.49 | 0.50 |
| Audio Flamingo 3 (think mode) | 0.50 | 0.50 | 0.51 | 0.50 | 0.49 | 0.52 | 0.46 | 0.52 | 0.50 | 0.50 | 0.50 | 0.44 | 0.50 |
| R1-AQA (think mode) | 0.51 | 0.51 | 0.48 | 0.50 | 0.49 | 0.54 | 0.51 | 0.47 | 0.52 | 0.60 | 0.56 | 0.50 | 0.52 |
| Step-Audio 2 mini Think | 0.48 | 0.44 | 0.51 | 0.53 | 0.53 | 0.49 | 0.51 | 0.46 | 0.50 | 0.70 | 0.67 | 0.46 | 0.52 |
| Audio-Reasoner | 0.65 | 0.62 | 0.55 | 0.57 | 0.56 | 0.49 | 0.49 | 0.49 | 0.48 | 0.65 | 0.54 | 0.64 | 0.56 |
| Step-Audio-R1 | 0.70 | 0.65 | 0.54 | 0.50 | 0.50 | 0.55 | 0.44 | 0.67 | 0.57 | 0.96 | 0.58 | 0.80 | 0.62 |
| MiMo-Audio (think mode) | 0.79 | 0.74 | 0.64 | 0.46 | 0.67 | 0.54 | 0.50 | 0.57 | 0.61 | 0.69 | 0.55 | 0.74 | 0.63 |
| Omni Language Models (OLMs) | | | | | | | | | | | | | |
| OpenOmni | 0.43 | 0.48 | 0.51 | 0.50 | 0.51 | 0.52 | 0.52 | 0.49 | 0.52 | 0.50 | 0.54 | 0.50 | 0.50 |
| Baichuan-Omni-1.5 | 0.56 | 0.47 | 0.54 | 0.52 | 0.48 | 0.51 | 0.51 | 0.44 | 0.51 | 0.57 | 0.46 | 0.52 | 0.51 |
| VITA-1.5 | 0.53 | 0.52 | 0.54 | 0.52 | 0.52 | 0.51 | 0.51 | 0.50 | 0.49 | 0.48 | 0.44 | 0.51 | 0.51 |
| HumanOmni | 0.50 | 0.52 | 0.54 | 0.51 | 0.47 | 0.48 | 0.52 | 0.48 | 0.49 | 0.80 | 0.50 | 0.50 | 0.53 |
| Ming-Lite-Omni-1.5 | 0.56 | 0.53 | 0.55 | 0.50 | 0.61 | 0.60 | 0.50 | 0.48 | 0.50 | 0.92 | 0.56 | 0.69 | 0.58 |
| Ola | 0.72 | 0.73 | 0.68 | 0.49 | 0.55 | 0.55 | 0.54 | 0.48 | 0.47 | 0.97 | 0.47 | 0.58 | 0.60 |
| Qwen2.5-Omni | 0.60 | 0.89 | 0.73 | 0.54 | 0.54 | 0.50 | 0.50 | 0.52 | 0.42 | 1.00 | 0.69 | 0.55 | 0.62 |
| Qwen3-Omni-Instruct | 0.83 | 0.92 | 0.80 | 0.64 | 0.73 | 0.66 | 0.52 | 0.60 | 0.52 | 0.99 | 0.70 | 0.65 | 0.71 |
| Gemini-2.5-Flash | 0.77 | 0.82 | 0.68 | 0.60 | 0.54 | 0.64 | 0.50 | 0.54 | 0.75 | 0.91 | 0.61 | 0.78 | 0.68 |

%abstention

Table 13: Comparison task accuracy across audio attributes.

| Models | Spectral & Amplitude | | | | Temporal | | Spatial& Environment | | | Timbre | | Scene Level | Avg. |
|--------------------------------------|----------------------|------------|----------|----------|----------|-------|----------------------|----------|---------------|---------|--------|-------------|------|
| | Pitch | Brightness | Loudness | Velocity | Duration | Tempo | Direction | Distance | Reverberation | Texture | Timbre | Counting | |
| Random Guess | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 |
| Human | 0.87 | 0.93 | 0.77 | 0.83 | 0.83 | 0.77 | 0.83 | 0.73 | 1.00 | 1.00 | 0.93 | 1.00 | 0.88 |
| Large Audio Language Models (LALMs) | | | | | | | | | | | | | |
| BAT | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.03 | 0.11 | 0.06 | 0.02 | 0.00 | 0.02 |
| MU-LLaMA | 0.00 | 0.03 | 0.00 | 0.06 | 0.00 | 0.11 | 0.00 | 0.01 | 0.10 | 0.39 | 0.48 | 0.01 | 0.10 |
| LLaMA-Omni | 0.43 | 0.49 | 0.45 | 0.47 | 0.46 | 0.48 | 0.49 | 0.40 | 0.40 | 0.54 | 0.50 | 0.51 | 0.47 |
| GLM-4-Voice | 0.50 | 0.50 | 0.51 | 0.55 | 0.48 | 0.50 | 0.50 | 0.44 | 0.46 | 0.49 | 0.50 | 0.38 | 0.48 |
| Audio Flamingo 3 | 0.46 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.49 | 0.48 | 0.50 | 0.50 | 0.50 | 0.41 | 0.49 |
| Voxtral-Mini | 0.52 | 0.43 | 0.49 | 0.46 | 0.51 | 0.51 | 0.50 | 0.46 | 0.49 | 0.55 | 0.49 | 0.50 | 0.49 |
| LLaMA-Omni2 | 0.54 | 0.50 | 0.50 | 0.44 | 0.50 | 0.44 | 0.50 | 0.52 | 0.50 | 0.50 | 0.47 | 0.46 | 0.49 |
| VITA-Audio-Plus-Vanilla | 0.50 | 0.50 | 0.50 | 0.51 | 0.50 | 0.51 | 0.47 | 0.53 | 0.50 | 0.54 | 0.48 | 0.40 | 0.50 |
| SALMONN | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.44 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 |
| Audio Flamingo 2 | 0.50 | 0.54 | 0.50 | 0.50 | 0.48 | 0.50 | 0.50 | 0.44 | 0.50 | 0.60 | 0.56 | 0.50 | 0.51 |
| Qwen2-Audio-Instruct | 0.52 | 0.49 | 0.47 | 0.54 | 0.50 | 0.51 | 0.48 | 0.52 | 0.50 | 0.67 | 0.60 | 0.47 | 0.52 |
| Baichuan-Audio-Instruct | 0.49 | 0.50 | 0.53 | 0.50 | 0.50 | 0.49 | 0.42 | 0.39 | 0.55 | 0.81 | 0.54 | 0.49 | 0.52 |
| Step-Audio 2 mini | 0.57 | 0.62 | 0.41 | 0.53 | 0.54 | 0.47 | 0.43 | 0.46 | 0.52 | 0.98 | 0.65 | 0.52 | 0.56 |
| R1-AQA | 0.64 | 0.72 | 0.57 | 0.52 | 0.50 | 0.47 | 0.49 | 0.51 | 0.54 | 0.99 | 0.79 | 0.46 | 0.60 |
| MiMo-Audio-Instruct | 0.81 | 0.72 | 0.58 | 0.51 | 0.56 | 0.50 | 0.50 | 0.48 | 0.59 | 0.93 | 0.58 | 0.49 | 0.60 |
| MiDashengLM | 0.65 | 0.75 | 0.72 | 0.46 | 0.68 | 0.55 | 0.44 | 0.48 | 0.50 | 0.96 | 0.69 | 0.44 | 0.61 |
| Kimi-Audio-Instruct | 0.87 | 0.81 | 0.70 | 0.50 | 0.75 | 0.55 | 0.44 | 0.48 | 0.95 | 0.99 | 0.68 | 0.69 | 0.70 |
| GPT-4o-Audio | 0.71 | 0.78 | 0.49 | 0.52 | 0.54 | 0.51 | 0.51 | 0.54 | 0.50 | 0.92 | 0.49 | 0.68 | 0.60 |
| Large Audio Reasoning Models (LARMs) | | | | | | | | | | | | | |
| Mellow | 0.11 | 0.19 | 0.05 | 0.08 | 0.24 | 0.20 | 0.12 | 0.15 | 0.28 | 0.32 | 0.39 | 0.38 | 0.21 |
| GAMA | 0.21 | 0.20 | 0.21 | 0.16 | 0.30 | 0.28 | 0.31 | 0.19 | 0.39 | 0.76 | 0.46 | 0.30 | 0.31 |
| Audio Flamingo 2 Sound-CoT | 0.26 | 0.53 | 0.34 | 0.22 | 0.52 | 0.48 | 0.50 | 0.45 | 0.50 | 0.66 | 0.47 | 0.49 | 0.45 |
| Audio Flamingo 3 (think mode) | 0.52 | 0.47 | 0.50 | 0.50 | 0.48 | 0.44 | 0.51 | 0.44 | 0.50 | 0.50 | 0.51 | 0.45 | 0.49 |
| R1-AQA (think mode) | 0.65 | 0.45 | 0.45 | 0.51 | 0.50 | 0.48 | 0.52 | 0.47 | 0.51 | 0.94 | 0.66 | 0.24 | 0.53 |
| Audio-Reasoner | 0.69 | 0.84 | 0.57 | 0.51 | 0.53 | 0.48 | 0.46 | 0.43 | 0.52 | 0.78 | 0.53 | 0.17 | 0.54 |
| Step-Audio 2 mini Think | 0.52 | 0.62 | 0.52 | 0.41 | 0.53 | 0.54 | 0.49 | 0.44 | 0.52 | 0.92 | 0.46 | 0.54 | 0.54 |
| Step-Audio-R1 | 0.57 | 0.68 | 0.54 | 0.50 | 0.50 | 0.58 | 0.51 | 0.53 | 0.83 | 0.95 | 0.62 | 0.73 | 0.63 |
| MiMo-Audio (think mode) | 0.83 | 0.76 | 0.56 | 0.49 | 0.71 | 0.51 | 0.49 | 0.48 | 0.73 | 0.91 | 0.68 | 0.71 | 0.66 |
| Omni Language Models (OLMs) | | | | | | | | | | | | | |
| VITA-1.5 | 0.49 | 0.49 | 0.50 | 0.45 | 0.50 | 0.51 | 0.46 | 0.41 | 0.61 | 0.47 | 0.43 | 0.41 | 0.48 |
| OpenOmni | 0.50 | 0.48 | 0.50 | 0.51 | 0.53 | 0.48 | 0.49 | 0.44 | 0.53 | 0.51 | 0.53 | 0.51 | 0.50 |
| Baichuan-Omni-1.5 | 0.58 | 0.47 | 0.51 | 0.51 | 0.48 | 0.49 | 0.49 | 0.39 | 0.46 | 0.84 | 0.56 | 0.59 | 0.53 |
| Ming-Lite-Omni-1.5 | 0.57 | 0.50 | 0.50 | 0.50 | 0.44 | 0.50 | 0.51 | 0.45 | 0.56 | 0.95 | 0.57 | 0.61 | 0.56 |
| HumanOmni | 0.59 | 0.70 | 0.50 | 0.46 | 0.53 | 0.52 | 0.52 | 0.36 | 0.49 | 0.91 | 0.64 | 0.68 | 0.58 |
| Ola | 0.75 | 0.72 | 0.60 | 0.46 | 0.57 | 0.48 | 0.43 | 0.45 | 0.68 | 0.97 | 0.59 | 0.58 | 0.61 |
| Qwen2.5-Omni | 0.82 | 0.84 | 0.76 | 0.52 | 0.60 | 0.53 | 0.49 | 0.52 | 0.69 | 0.97 | 0.80 | 0.60 | 0.68 |
| Qwen3-Omni-Instruct | 0.90 | 0.83 | 0.70 | 0.49 | 0.80 | 0.64 | 0.49 | 0.42 | 0.83 | 0.98 | 0.75 | 0.99 | 0.74 |
| Gemini-2.5-Flash | 0.77 | 0.75 | 0.57 | 0.54 | 0.81 | 0.55 | 0.47 | 0.46 | 0.68 | 0.97 | 0.63 | 0.63 | 0.65 |

%abstention

Table 14: Recognition task accuracy across audio attributes.