

# DYNAMIC-SUPERB PHASE-2: A COLLABORATIVELY EXPANDING BENCHMARK FOR MEASURING THE CAPABILITIES OF SPOKEN LANGUAGE MODELS WITH 180 TASKS

Chien-yu Huang<sup>1</sup>, Wei-Chih Chen<sup>1</sup>, Shu-wen Yang<sup>1</sup>, Andy T. Liu<sup>1</sup>, Chen-An Li<sup>1</sup>, Yu-Xiang Lin<sup>1</sup>, Wei-Cheng Tseng<sup>2</sup>, Anuj Diwan<sup>2</sup>, Yi-Jen Shih<sup>2</sup>, Jiatong Shi<sup>3</sup>, William Chen<sup>3</sup>, Chih-Kai Yang<sup>1</sup>, Xuanjun Chen<sup>1</sup>, Chi-Yuan Hsiao<sup>1</sup>, Puyuan Peng<sup>2</sup>, Shih-Heng Wang<sup>1</sup>, Chun-Yi Kuan<sup>1</sup>, Ke-Han Lu<sup>1</sup>, Kai-Wei Chang<sup>1</sup>, Fabian Ritter-Gutierrez<sup>4</sup>, Kuan-Po Huang<sup>1</sup>, Siddhant Arora<sup>3</sup>, You-Kuan Lin<sup>1</sup>, Ming To Chuang<sup>1</sup>, Eunjung Yeo<sup>3</sup>, Calvin Chang<sup>3</sup>, Chung-Ming Chien<sup>5</sup>, Kwanghee Choi<sup>3</sup>, Cheng-Hsiu Hsieh<sup>1</sup>, Yi-Cheng Lin<sup>1</sup>, Chee-En Yu<sup>1</sup>, I-Hsiang Chiu<sup>1</sup>, Heitor R. Guimarães<sup>6</sup>, Jionghao Han<sup>3</sup>, Tzu-Quan Lin<sup>1</sup>, Tzu-Yuan Lin<sup>7</sup>, Homu Chang<sup>1</sup>, Ting-Wu Chang<sup>1</sup>, Chun Wei Chen<sup>1</sup>, Shou-Jen Chen<sup>1</sup>, Yu-Hua Chen<sup>1</sup>, Hsi-Chun Cheng<sup>1</sup>, Kunal Dhawan<sup>8</sup>, Jia-Lin Fang<sup>1</sup>, Shi-Xin Fang<sup>1</sup>, Kuan-Yu Fang Chiang<sup>1</sup>, Chi An Fu<sup>1</sup>, Hsien-Fu Hsiao<sup>1</sup>, Ching Yu Hsu<sup>1</sup>, Shao-Syuan Huang<sup>1</sup>, Lee Chen Wei<sup>1</sup>, Hsi-Che Lin<sup>1</sup>, Hsuan-Hao Lin<sup>1</sup>, Hsuan-Ting Lin<sup>1</sup>, Jian-Ren Lin<sup>1</sup>, Ting-Chun Liu<sup>1</sup>, Li-Chun Lu<sup>1</sup>, Tsung-Min Pai<sup>1</sup>, Ankita Pasad<sup>5</sup>, Shih-Yun Shan Kuan<sup>1</sup>, Suwon Shon<sup>9</sup>, Yuxun Tang<sup>10</sup>, Yun-Shao Tsai<sup>1</sup>, Jui-Chiang Wei<sup>1</sup>, Tzu-Chieh Wei<sup>1</sup>, Chengxi Wu<sup>1</sup>, Dien-Ruei Wu<sup>1</sup>, Chao-Han Huck Yang<sup>8</sup>, Chieh-Chi Yang<sup>1</sup>, Jia Qi Yip<sup>4</sup>, Shao-Xiang Yuan<sup>1</sup>, Haibin Wu<sup>1</sup>, Karen Livescu<sup>5</sup>, David Harwath<sup>2</sup>, Shinji Watanabe<sup>3</sup>, Hung-yi Lee<sup>1</sup>

<sup>1</sup>National Taiwan University   <sup>2</sup>University of Texas at Austin   <sup>3</sup>Carnegie Mellon University

<sup>4</sup>Nanyang Technological University   <sup>5</sup>Toyota Technological Institute at Chicago

<sup>6</sup>Université du Québec (INRS-EMT)   <sup>7</sup>Independent Researcher   <sup>8</sup>NVIDIA   <sup>9</sup>ASAPP

<sup>10</sup>Renmin University of China

## ABSTRACT

Multimodal foundation models, such as Gemini and ChatGPT, have revolutionized human-machine interactions by seamlessly integrating various forms of data. Developing a universal spoken language model that comprehends a wide range of natural language instructions is critical for bridging communication gaps and facilitating more intuitive interactions. However, the absence of a comprehensive evaluation benchmark poses a significant challenge. We present Dynamic-SUPERB *Phase-2*, an open and evolving benchmark for the comprehensive evaluation of instruction-based universal speech models. Building upon the first generation, this second version incorporates 125 new tasks contributed collaboratively by the global research community, expanding the benchmark to a total of *180 tasks*, making it the largest benchmark for speech and audio evaluation. While the first generation of Dynamic-SUPERB was limited to classification tasks, Dynamic-SUPERB *Phase-2* broadens its evaluation capabilities by introducing a wide array of novel and diverse tasks, including regression and sequence generation, across speech, music, and environmental audio. Evaluation results show that no model performed well universally. SALMONN-13B excelled in English ASR and Qwen2-Audio-7B-Instruct showed high accuracy in emotion recognition, but current models still require further innovations to handle a broader range of tasks. We open-source the project at <https://github.com/dynamic-superb/dynamic-superb>.

## 1 INTRODUCTION

Recent advancements in large language models (LLMs) have accelerated the development of natural language processing (NLP) (Touvron et al., 2023a; Achiam et al., 2023; Li et al., 2023b; Anthropic,

2023; Bai et al., 2023). These models can follow natural language instructions, making users quickly adopt them for a variety of applications. They have been integrated into commercial products, such as ChatGPT (Achiam et al., 2023) and Claude (Anthropic, 2023), as well as in the open-source research community, including the LLaMA series (Touvron et al., 2023a;b; Dubey et al., 2024). Yet, they are primarily text-based models, meaning they cannot process speech or audio, which are essential for more natural ways of communication and interaction with the real world.

Compared to written text, spoken language has always been a more natural and convenient way for humans to communicate. Spoken language conveys a wealth of information, including semantics, prosody, emotion, and speaker characteristics, while text is limited to representing semantic information, which can sometimes even depend on the prosodic cues present in spoken language (Lin et al., 2024). This highlights the need for universal speech models and explains why automatic speech recognition (ASR) systems using text-based language models are not optimal. Several attempts have been made to develop instruction-based universal speech or audio models capable of performing various tasks, such as LTU-AS (Gong et al., 2023), SALMONN (Tang et al., 2024a), Qwen-Audio (Chu et al., 2023; 2024), and WavLLM (Hu et al., 2024).

Despite significant research in universal speech models, evaluating them effectively and comprehensively remains a major challenge. In NLP, benchmarks for text LLMs include a large number of tasks. CrossFit (Ye et al., 2021) includes 160 tasks, BIG-bench (Srivastava et al., 2022) comprises 204 tasks, and Natural-Instructions (Mishra et al., 2022; Wang et al., 2022) provides 1,616 tasks. For universal speech models, several benchmarks have been proposed but with a *fixed* and *limited* set of tasks (typically around a dozen) to evaluate specific capabilities of models (Yang et al., 2021; Dunbar et al., 2022; Maimon et al., 2024). This is insufficient for evaluating a universal model, as we aim to investigate whether models can handle a broader range of tasks beyond fixed benchmarks. There is a strong demand for a benchmark that can evaluate these models across various aspects with a wide range of speech tasks, which led to the creation of Dynamic-SUPERB (Huang et al., 2024a).

Dynamic-SUPERB is the first benchmark for evaluating universal instruction-based speech and audio models. As models advance, we *dynamically* expand the benchmark by adding new tasks to provide better guidance for researchers in model development. We have designed a well-structured pipeline that harnesses the collective efforts of the community to gather diverse challenging, novel, and creative tasks. The first generation of Dynamic-SUPERB consists of 55 tasks, covering various aspects of speech (e.g., semantics, speakers, etc.), making it the speech benchmark with the most tasks. However, this is not sufficient to pave the way for developing universal speech models, especially given the complexity and richness of the information conveyed by spoken language.

This paper presents the **Dynamic-SUPERB Phase-2**. We collected *125 new tasks* with contributions from the global research community, allowing the benchmark to grow more than twice as large as its first iteration. They cover a broad range of types, and some tasks present novel challenges that have not been explored in any previous research. Besides, previous speech research has typically treated music and environmental audio as background noise to be ignored. However, these sounds contain rich information and share overlapping elements that complement each other. Thus, in Dynamic-SUPERB, we also introduced preliminary music and audio tasks from established music and audio benchmarks (Yuan et al., 2023; Turian et al., 2022). We define the core tasks by reformulating the SUPERB (Yang et al., 2021) (speech), MARBLE (Yuan et al., 2023) (music), and HEAR (Turian et al., 2022) (audio) benchmarks for quick-round research experiments. One challenge of evaluating with such a large-scale benchmark is deriving concrete and useful insights from hundreds of evaluation results. We provide a taxonomy for every task in Dynamic-SUPERB, where tasks are clustered by the specific model capabilities they probe. Researchers can follow this taxonomy to develop or reinforce specific capabilities of the models they build. To our knowledge, Dynamic-SUPERB is the largest benchmark and the first to provide such a detailed task taxonomy in speech processing.

We conducted a comprehensive evaluation of several models using Dynamic-SUPERB *Phase-2*. To handle the diverse output formats of these models, we propose an automated LLM-based pipeline for general evaluation across tasks. The evaluation results show that these models perform well only on a limited range of tasks. For example, SALMONN-13B performed well on English ASR, while WavLLM achieved high accuracy in emotion recognition. However, they do not generalize well to a much broader range of tasks in speech recognition and paralinguistics. Notably, training on diverse data, even with significant differences in signal-level characteristics, can enhance performance across domains. We observed that spoken language models outperformed music models in certain

music tasks. This highlights the potential for developing unified models for speech, music, and general audio. We have open-sourced all materials to support reproducibility and further research. We invite researchers to help expand Dynamic-SUPERB, making it more diverse and comprehensive to advance universal spoken language models.

## 2 RELATED WORKS

### 2.1 INSTRUCTION-FOLLOWING UNIVERSAL SPEECH MODELS

LLMs have shown strong natural language processing abilities and are used in speech, audio, and music applications. Recent frameworks integrate a pre-trained speech encoder with an LLM through fine-tuning techniques such as LoRA (Hu et al., 2021). LTU-AS (Gong et al., 2023) combines Whisper (Radford et al., 2023) with LLaMA (Touvron et al., 2023a), and it is fine-tuned on open-ended speech and audio question-answering. SALMONN adopts a window-level Q-former (Li et al., 2023a) to generate soft prompts fusing representations from the speech and audio encoders. Qwen-Audio (Chu et al., 2023) introduces task-specific tags into Qwen to encourage knowledge sharing and prevent interference across diverse tasks, and it supports multi-turn dialogues for both audio and text inputs. WavLLM (Hu et al., 2024) uses curriculum learning to prevent overfitting on specific tasks while maintaining the LLM’s original capabilities. All these models accept speech, audio, and text as input but output only text, focusing on *understanding* rather than *generation* in speech and audio. Several attempts have also used prompts for speech, music, and audio generation (Guo et al., 2023; Agostinelli et al., 2023; Liu et al., 2023a). However, to our knowledge, there is no universal model capable of handling both generation and understanding tasks. Moreover, these models have not been comprehensively evaluated on a benchmark, which hinders fair comparisons between them.

### 2.2 EVALUATION BENCHMARKS

Several benchmarks have been developed to evaluate speech models. SUPERB (Yang et al., 2021) is the most widely used benchmark for assessing the performance of speech foundation models across various tasks that cover different aspects of speech. On the other hand, SLUE (Shon et al., 2022; 2023) focuses more specifically on spoken language understanding. Yet, they are primarily limited to English. LeBenchmark (Evain et al., 2021; Parcollet et al., 2023) evaluates speech foundation models specifically for French, while IndicSUPERB (Javed et al., 2023) is dedicated to Indian languages. ML-SUPERB (Shi et al., 2023; 2024b) offers ASR and language identification tasks in 100+ languages. Aside from speech, MARBLE (Yuan et al., 2023) and HEAR (Turian et al., 2022) offer platforms for evaluating various music and a wide range of audio tasks. Using these benchmarks, researchers build a specialized model for each task. Thus, the number of tasks is limited due to the growing costs associated with adding more tasks. Conversely, universal models are expected to do various tasks without fine-tuning for each one, allowing users to engage in far more diverse applications and enabling benchmark developers to include more tasks for comprehensive evaluation. In NLP and CV, benchmarks have been developed to evaluate models across a much broader range of tasks (Bitton et al., 2023; Srivastava et al., 2022; Mishra et al., 2022). However, in speech, we lack a benchmark of comparable scale. Dynamic-SUPERB (Huang et al., 2024a) is the first benchmark for instruction-based universal speech models. While its first version includes far more tasks than all other speech benchmarks, they are all classification tasks. AIR-bench (Yang et al., 2024b), on the other hand, includes tasks from speech, music, and audio, and expands beyond classification, although the number of tasks available for evaluating universal models remains limited (19 tasks in foundation track). Thus, there remains a critical need for a comprehensive benchmark that evaluates universal models across a broader range of tasks to fully assess their capabilities. Consequently, we developed Dynamic-SUPERB *Phase-2*, substantially upgrading the first generation with a detailed taxonomy and establishing community contribution protocols to facilitate continual task integration.

## 3 DYNAMIC-SUPERB PHASE-2

### 3.1 OVERVIEW

Figure 1 depicts the framework of Dynamic-SUPERB *Phase-2*. Our goal is to evaluate universal models that meet the following criteria: (1) The model can accept speech, music, or other audio as

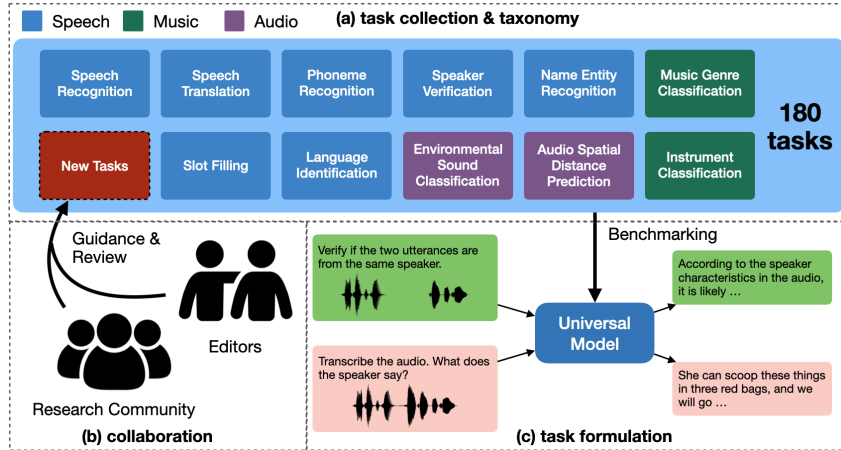


Figure 1: An overview of Dynamic-SUPERB.

input. (2) The model can perform specific tasks without requiring fine-tuning. (3) The model follows natural language instructions to execute the corresponding tasks. To comprehensively evaluate these universal models, we have collected hundreds of tasks and developed a task taxonomy to better guide benchmark users (Figure 1a). All tasks in Dynamic-SUPERB *Phase-2* are intended solely for testing purposes; we do not provide training data for two main reasons. First, current models are trained on large-scale, open-source, or proprietary datasets, making it challenging to ensure that all models are trained under consistent conditions. Second, we aim to evaluate universal models that do not require fine-tuning for downstream tasks.

The Dynamic-SUPERB project is designed to evolve dynamically with research advancements by incorporating novel tasks from the research community (Figure 1b). With the first call for tasks, Dynamic-SUPERB *Phase-2* has grown from its first generation, expanding from 55 to 180 tasks. Table 1 compares several benchmarks used in speech, music, and audio research, demonstrating that Dynamic-SUPERB *Phase-2* is the largest benchmark covering all three areas. It provides a more fine-grained evaluation than any other benchmark.

Table 1: Comparison of popular benchmarks with Dynamic-SUPERB.

| Benchmark  | SUPERB | SLUE | HEAR | MARBLE | AIR-Bench | Dynamic-SUPERB Phase-1 | Dynamic-SUPERB Phase-2 |
|------------|--------|------|------|--------|-----------|------------------------|------------------------|
| # of tasks | 13     | 7    | 19   | 13     | 19*       | 55                     | 180                    |
| Speech     | ✓      | ✓    |      |        | ✓         | ✓                      | ✓                      |
| Music      |        |      |      | ✓      | ✓         | ✓                      | ✓                      |
| Audio      |        |      | ✓    |        | ✓         | ✓                      | ✓                      |

\* foundation track

### 3.2 TASK FORMULATION

As Figure 1c depicts, in Dynamic-SUPERB *Phase-2*, each task is structured to include: (1) Text instruction: A natural language instruction that guides the model on the task to perform. Each task has multiple different instructions to evaluate the model’s ability to understand instructions. (2) Audio component: At least one audio element, which can be in the input, the output, or both. (3) (Optional) Text component: Text elements (other than instruction) that may serve as inputs or outputs. The number of audio elements in the inputs or outputs may vary depending on the task. For example, in speaker verification, two utterances are used to determine whether they were produced by the same speaker (green blocks in Figure 1c). Besides, text inputs are not always required; for instance, ASR does not involve any text inputs. We use text format for instructions instead of spoken format. Spoken instructions involve varying content, speaker characteristics, and prosody, making them more complex. Text instructions act as an intermediary, bridging text and spoken universal models. These designs ensure consistency and simplify the model’s understanding and processing of diverse tasks while maintaining the benchmark’s extensibility with minimal constraints.

Another problem in evaluating universal models is the format of classification and regression tasks. Generally, a task-specialized model generates predicted labels (soft distributions or hard labels) within a pre-defined set of labels for classification or numeric values for regression within a pre-defined range. However, a universal model does not necessarily have access to this information, and thus it is hard to evaluate them in the same way, especially using it to perform several different tasks. To address this, we believe that a universal speech model should produce outputs in natural language for all tasks (model outputs in Figure 1c). Therefore, in Dynamic-SUPERB, *all classification labels are represented as text*. For regression tasks with specific formats (e.g., scalars, JSON, or Python-style lists), we can parse the natural language outputs using a post-processing pipeline (Section 4.2).

### 3.3 CALL FOR TASKS

Dynamic-SUPERB fosters community collaboration, encouraging the addition of innovative tasks to remain relevant in this rapidly changing field. In *Phase-2*, we initiated a call for tasks in March 2024 to invite contributions from the research community. We established an organized and transparent submission process on our GitHub portal, where we organize members who serve as editors to guide contributors. Contributors propose tasks by providing relevant information on GitHub. Each proposal is assigned to an editor who checks for major issues and prevents duplicated efforts. Then, task proposers upload their data to our Huggingface space in specified formats, using only datasets with proper licenses. They complete submissions by opening a pull request on GitHub with the necessary files. Afterward, editors review submissions on a rolling basis, offering suggestions for refinement rather than immediate acceptance or rejection. After iterative improvements, accepted tasks are merged into the repository. Between March and July 2024, we received about 145 task proposals and accepted 91 tasks. Tasks still under review are not included here but will be featured in the next phase. For details, please refer to Appendix F.

### 3.4 TASK TAXONOMY

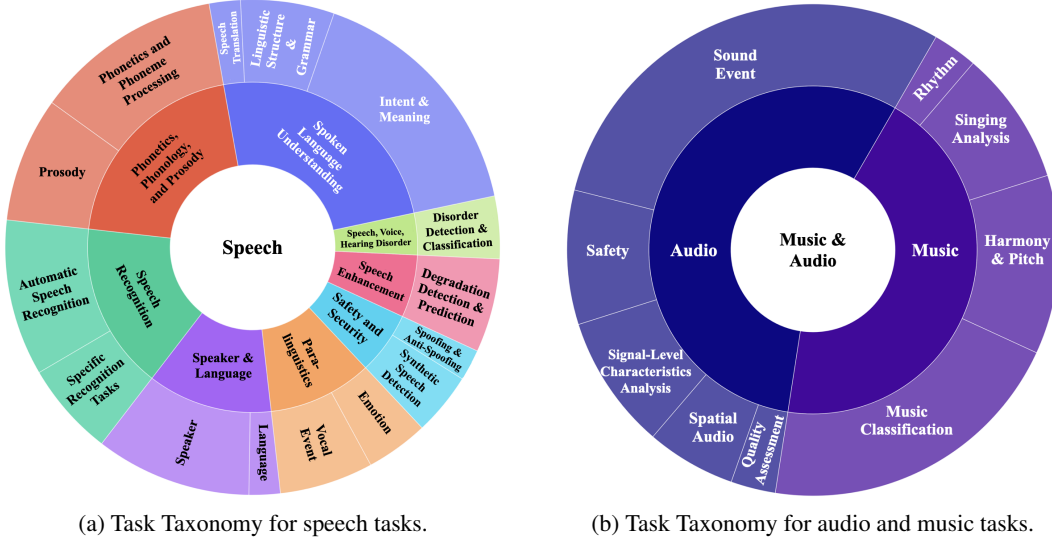


Figure 2: Task taxonomy in Dynamic-SUPERB.

One primary challenge in building a large-scale benchmark is offering valuable insights to its users. To address this, we developed a task taxonomy<sup>1</sup> that helps researchers interpret performance results across various tasks. Researchers can leverage this taxonomy to select specific tasks for model development instead of evaluating every task in the benchmark.

Figure 2 shows the high-level task taxonomy in Dynamic-SUPERB *Phase-2*. Due to the space limitation, we only show some representative tasks. Each leaf node includes at least one task and

<sup>1</sup>We developed the taxonomy by referencing sessions from the INTERSPEECH conference and EDICS of IEEE SPS. Though not identical, they have a very similar structure (please refer to Appendix A for comparison).

may be further categorized into more fine-grained sub-domains, which are not shown here due to space constraints. For example, within ‘Speaker & Language/Speaker’, we have two categories: ‘Speaker Characteristics’ and ‘Speaker Identification’, each containing several tasks. Please refer to Appendix B for the complete task list. We first categorize tasks into two primary fields: (I) speech and (II) music & audio, which are generally distinguished by the source of the sound. For instance, speech is produced by human vocal cords, music is created using instruments, and audio includes sounds generated by other creatures, materials, or natural phenomena. We then split each field into several domains based on the key attributes and challenges that the tasks within them present.

### 3.4.1 SPEECH

As Figure 2a shows, there are 8 domains within the speech field. **Speech Recognition** focuses on converting spoken language into text. This includes tackling various challenges like multilingual and code-switching ASR, as well as spontaneous ASR which is very different from audiobook-style ones. It also contains specialized tasks such as command recognition and keyword spotting. **Speaker and Language** addresses the analysis of speaker characteristics and languages, covering tasks such as speaker verification, diarization, and language identification. **Spoken Language Understanding** deals with understanding and analysis of the content and semantics of spoken language. It covers tasks like sentiment analysis and speech translation. **Phonetics, Phonology, and Prosody** focuses on the sound structure of speech, including phoneme recognition, pronunciation evaluation, and prosodic features like stress and accent classification. **Paralinguistics** explores non-verbal aspects of speech, such as emotion recognition and vocal event detection, which can capture nuances like screaming or coughing. **Speech Enhancement** aims to improve speech quality by detecting and mitigating noise, reverberation, and other degradations, but currently, we only have understanding tasks. **Speech, Voice, and Hearing Disorders** is dedicated to identifying and classifying disorders such as stuttering and so on. **Safety and Security** focuses on detecting synthetic or manipulated speech, addressing tasks like spoof detection and recognizing deepfake voices.

### 3.4.2 AUDIO & MUSIC

The audio and music field includes a wide range of tasks that focus on various attributes of sound beyond speech, such as musical elements, environmental sounds, and advanced sound analysis. This field is divided into 9 domains, each addressing a specific aspect of audio or music processing. **Music Classification** tasks focus on categorizing musical elements such as instruments, genres, and emotions, providing a foundation for recognizing and analyzing different types of musical content. **Pitch Analysis** delves into identifying the pitch and harmony within music, including tasks like pitch extraction, chord classification, and key detection. **Rhythm Analysis** involves tasks such as beat tracking, which is critical for understanding the temporal structure of music. In **Singing Analysis**, tasks address both lyric recognition and translation, as well as the classification of vocal techniques used in singing. **Quality Assessment** evaluates the perceived quality of singing, including automated predictions of Mean Opinion Scores (MOS). The **Sound Event** domain is broader, covering various sound sources such as animals, the environment, and human activities. Tasks range from animal sound classification to emergency traffic detection and even advanced tasks like multi-channel sound event understanding. **Safety** domain includes detecting deepfakes in singing voices and identifying manipulated audio files, ensuring sound authenticity and integrity. **Spatial Audio** covers all tasks related to understanding spatial information, such as estimating the distance or position of sounds in the real world. Finally, **Signal Characteristics Analysis** domain addresses general signal-level characteristics of audio, including sound effect detection, and duration prediction.

### 3.4.3 CORE TASKS

While our task taxonomy provides a comprehensive, hierarchical, and systematic approach for benchmark users to easily get started with Dynamic-SUPERB *Phase-2*, evaluating all tasks with limited resources remains a challenge. To address this, alongside the task taxonomy, we have selected **core tasks** for speech, music, and audio. These core tasks are reduced subsets of essential tasks that have been widely used or studied within the research community. We include tasks from three popular benchmarks: SUPERB (speech), MARBLE (music), and HEAR (audio). The three benchmarks were originally designed for evaluating encoders, not for an instruction-following universal model, so modifications are required. We crafted instructions for each task, reduced the data

size, and reformulated them into the Dynamic-SUPERB format. Besides, we replaced the datasets in some tasks to address licensing issues. Since the core task set is much smaller than the full benchmark, researchers can more efficiently evaluate their models across a reasonable range of domains.

## 4 EXPERIMENTAL SETTINGS

### 4.1 MODELS

We evaluated several publicly-available models on Dynamic-SUPERB *Phase-2*: SALMONN, LTU-AS, Qwen-Audio-Chat, Qwen2-Audio-7B-Instruct, WavLLM, MU-LLaMA (Liu et al., 2024), GAMA (Deshmukh et al., 2024). SALMONN is further categorized into 7B and 13B versions based on the size of its LLM component. All of these models are publicly available, and we utilized their official implementations for inference without any modifications. The first five models are instruction-based speech models, each also trained with audio understanding capabilities. However, as they were not trained on music data, we do not expect them to perform well on music-related tasks. Hence, we further included MU-LLaMA and GAMA. MU-LLaMA is specifically designed for various music-understanding tasks, and GAMA is trained with both general audio data and a music corpus. Similar to the speech models, MU-LLaMA and GAMA adopt an LLM-based framework. They use a music or audio encoder, such as MERT (LI et al., 2024) or Q-former, to convert music or audio into features, which the LLM then uses as prompts for subsequent reasoning. The sampling strategies used for generating outputs from the LLMs were retained as defined in their official implementations. Finally, we implemented a cascaded system baseline called Whisper-LLaMA. The system first transcribes audio using Whisper-v3-large, and then LLaMA3.1-8B processes the transcriptions to perform various tasks based on the provided instructions.

When evaluating models on these diverse tasks, several challenges inevitably arise. One challenge was testing the baseline models on tasks involving multiple audio inputs. For example, in speaker verification, a model determines whether two utterances are produced by the same speaker. Among the evaluated models, only Qwen-Audio and Qwen2-Audio provide interfaces that allow inputting multiple audio files. For the remaining models, we concatenated all audio files, separated by 0.5 seconds of silence, and used the concatenated file as input.

Last, different models have their own maximum supported audio durations. Models using Whisper to encode speech face an inherent limitation: Whisper can process audio files of up to 30 seconds, and some models are restricted to handling even shorter durations. Consequently, we retained the original model settings and did not modify the model architecture to accommodate longer audio clips. A preliminary analysis of all audio (Appendix G) in Dynamic-SUPERB *Phase-2* shows that only a small proportion (around 7%) exceeds 30 seconds in length. Hence, we believe our settings do not largely impact the evaluation results and ensure reasonable inference efficiency.

### 4.2 EVALUATION METRICS

The outputs of universal speech models are natural language sentences, making it difficult to assess their correctness using conventional metrics. For classification tasks, such as emotion recognition, a task-specific model outputs a label from predefined emotions in the dataset, while a universal model generates sentences like “The speaker sounds happy.” In this case, we can easily assess the correctness of the former by comparing labels, but this approach does not apply to the latter. Similarly, in regression tasks, a natural language response like “Using MOS scoring criteria, I give the audio a score of 3 out of 5” makes it difficult to directly use metrics such as mean square error.

In evaluation, we categorize tasks into three types: (1) classification, (2) regression, and (3) sequence generation. For each type, we use different pipelines to evaluate the models’ outputs. For classification tasks, we utilize external LLMs (GPT-4o) as *referees*, with the temperature set to be zero for evaluation consistency, to evaluate whether the outputs from speech or music models<sup>2</sup> match the ground truth. This approach has been widely adopted in the NLP community (Wang et al., 2023; Liu et al., 2023b; Chiang & Lee, 2023), and we extend its application to speech research. We design a prompt that includes the task instructions, the model’s output to be evaluated,

<sup>2</sup>To avoid ambiguity, unless specified otherwise, the terms “output” or “model output” in this section refer to the results generated by the models in Sec. 4.1.

and the corresponding ground truth. The LLM judge gets this prompt and determines if the output aligns with the ground truth using a chain-of-thought reasoning process. By leveraging the strong instruction-following capabilities of these LLMs, we constrain them to provide the final decision in a fixed format, allowing us to extract the answer using simple methods such as regular expressions<sup>3</sup>. We then define accuracy as the percentage of outputs that the LLM judge considers aligned with the ground truth. Please refer to Appendix E for details on the prompts and the alignment between LLMs and humans. **Importantly, although we used GPT-4o for evaluation in the main text, to support reproducibility we have also included comprehensive evaluation results using the open-source LLM (LLaMA3.1-8b-Instruct) in Appendix C.**

For regression tasks, the above method cannot be applied because performance is not assessed using hard labels. Instead, we use LLMs (GPT-4o) as a *post-processor* to transform natural language responses into a format compatible with the original metrics used in these tasks. Specifically, we developed a prompt that directs the LLMs to extract essential information from the output and convert it into the same format as the ground truth. If the output lacks correct or relevant information that can be converted to match the ground truth format, the LLM returns “N/A” to mark it as invalid. Since conventional metrics cannot accommodate “N/A” and setting a default value is unreasonable due to the variability of metrics, we also calculate the N/A rate for each regression task, defined as the percentage of invalid outputs within a task. A higher N/A rate indicates that the model struggles with following instructions, which may indirectly affect its performance on the task.

For sequence generation tasks such as ASR and captioning, we apply their original metrics directly to the raw outputs from baseline models. This is because identifying redundant prefixes or sections unrelated to the task objectives is highly challenging in sequence generation, even with human involvement. Besides, our review of the original evaluation results reported by these baseline models revealed no explicit mention of post-processing procedures. Therefore, we base our evaluations solely on the unprocessed outputs, ensuring consistency and objectivity across all models.

## 5 RESULTS



Figure 3: Domain-level relative-score-based performance comparison across different models. We exclude tasks where the models have a 100% NA rate or where Whisper-LLaMA scores zero, as these make it impossible to average their scores with other tasks.

<sup>3</sup>We have tried to ask the speech universal models to directly output their decisions in a fixed format, but they do not always follow the instructions. Therefore, the evaluation pipeline described here is necessary.



### 5.1 DOMAIN-LEVEL PERFORMANCE COMPARISON IN TAXONOMY

Here, we analyze the results based on the taxonomy. Due to space limitations, we can only report and compare the results of different models at the domain level in this subsection. Since different metrics are used across tasks, we adopted a **relative-score-based** method to summarize task performance. For each task, we calculated the relative improvement of each model compared to the cascaded system baseline (Whisper + LLaMA) and then obtained the domain-level scores by averaging all its improvements across the tasks within each domain. For regression tasks, we introduced the N/A rate to measure how well a model meets the task requirements. The original task metric was only computed on instances that followed the task format (after post-processing by the LLM). To account for whether a model can follow instructions, we incorporated the N/A rate into the reported scores. Accordingly, we calculated the improvement based on scaled values. For metrics where a higher value indicates better performance (such as the F1 score), we multiplied the metric value by  $(1 - \text{N/A rate})$ . Conversely, for metrics where a lower value indicates better performance (such as word error rate), we divided the metric value by  $(1 - \text{N/A rate})$ .

Figure 3 presents a domain-level comparison of relative scores across different models. Whisper-LLaMA consistently has zero scores throughout all domains, serving as the reference baseline for evaluating other models. In speech domains, no model outperforms Whisper-LLaMA in speech recognition and spoken language understanding. This shows that using ASR remains a strong baseline for language understanding, as text more explicitly represents semantic information than speech. However, Whisper-LLaMA lags behind SALMONN in phoneme recognition, a task for which SALMONN is explicitly trained, indicating that task-specific training still provides superior performance in specialized domains.

Conversely, in the speaker and paralinguistics domains, all models surpass the baseline. This improvement occurs because the ASR process in the cascaded system tends to discard critical information from the speech signal, such as speaker characteristics, pitch, and emotion. In contrast, other models utilize soft representations that retain more speech information, giving them the potential to learn beyond ASR, and some of these models were also trained to perform specific tasks within certain domains<sup>4</sup>. Additionally, music-specific large language models like GAMA-IT and Mu-LLaMA perform poorly on speech recognition and understanding tasks, which is expected since they were primarily designed for music understanding. Nevertheless, they surprisingly achieved scores comparable to the speech models in speaker and paralinguistics.

Turning to audio and music domains, where ASR models cannot transcribe non-speech information into text, most models outperform the baseline in several areas, such as music classification and sound event detection. Surprisingly, spoken language models primarily trained on speech data (Qwen models, SALMONN, WavLLM)<sup>5</sup> outperform the two music models (GAMA-IT and Mu-LLaMA) in various music-related domains such as harmony and pitch, music classification, and rhythm analysis. We speculate that training on diverse data, even with significant differences in signal-level characteristics, enhances performance across domains, emphasizing the importance of developing unified models for speech, music, and general audio processing.

Lastly, we observe some outliers with significantly higher or lower values in the figure, typically resulting from unbounded evaluation metrics. For instance, in the phonetics and prosody domain, WavLLM and LTU-AS exhibit poor scores due to their erroneous outputs in phone/phoneme segment counting tasks, predicting thousands of segments in utterances lasting only seconds and resulting in an excessively high mean square error. While the relative score-based approach effectively summarizes model performance, domain-level scores can be distorted by specific tasks within a domain. Thus, we encourage researchers aiming to develop model performance in specific domains to comprehensively report task-level scores to more accurately reflect their models’ capabilities.

### 5.2 CORE TASKS RESULTS

We tested the performance of the models on core tasks that are widely studied and commonly used. Table 2 presents the detailed results of the models on SUPERB tasks, displaying each task’s perfor-

<sup>4</sup>Although the models have been trained on some tasks in the same domains, the training tasks of these models generally do not use the same instructions as Dynamic-SUPERB and do not offer as broad coverage.

<sup>5</sup>Qwen and SALMONN are trained on a mixture of speech (dominant component), music, and audio data.

Table 2: Evaluation on SUPERB tasks. A “\*” indicates that the metric is scaled with NA rate, and a “-” means a model has a 100% NA rate.

|                      | PR    | KS    | IC    | ER    | ASR   | QbE   | SF    |         | SV    | SD      |
|----------------------|-------|-------|-------|-------|-------|-------|-------|---------|-------|---------|
|                      | PER ↓ | Acc ↑ | Acc ↑ | Acc ↑ | WER ↓ | Acc ↑ | F1* ↑ | CER* ↓  | Acc ↑ | DER* ↓  |
| Whisper-LLaMA        | 100.1 | 36.5  | 36.5  | 12.5  | 34.0  | 46.5  | 53.6  | 51.3    | 48.0  | 1068.7  |
| SALMONN-7B           | 25.4  | 30.5  | 21.0  | 12.1  | 15.1  | 49.0  | 61.6  | 61.5    | 45.0  | -       |
| SALMONN-13B          | 24.6  | 2.0   | 5.5   | 12.5  | 2.8   | 51.5  | 29.6  | 106.3   | 49.0  | -       |
| Qwen-Audio-Chat      | 100.6 | 60.5  | 26.5  | 70.4  | 69.4  | 48.0  | 42.2  | 90.1    | 51.0  | 8852.2  |
| Qwen2-Audio-7B-Inst. | 101.0 | 47.0  | 26.0  | 75.8  | 36.7  | 53.5  | 42.4  | 55.7    | 50.0  | 4664.4  |
| WavLLM               | 100.0 | 43.0  | 28.0  | 12.1  | 6.9   | 49.5  | 50.8  | 84.9    | 55.0  | 3621.2  |
| LTU-AS               | 102.8 | 1.0   | 0.1   | 25.8  | 96.3  | 45.0  | 10.2  | 559.3   | 47.0  | 14923.3 |
| GAMA-IT              | 100.4 | 2.0   | 1.0   | 0.0   | 116.7 | 2.5   | 2.7   | 4750.0  | 45.0  | 187.6   |
| Mu-LLaMA             | 110.3 | 4.0   | 3.5   | 12.9  | 103.7 | 51.0  | 1.0   | 15869.6 | 44.0  | -       |

**PR**: Phoneme Recognition, **KS**: Keyword Spotting, **IC**: Intent Classification, **ER**: Emotion Recognition, **QbE**: Query-by-Example, **SF**: Slot Filling, **SV**: Speaker Verification, **SD**: Speaker Diarization

mance in its own metric rather than relative scores. Please refer to Appendix D for results on HEAR and MARBLE. It is worth noting that these tasks have been adapted to the Dynamic-SUPERB framework. Therefore, the results are not directly comparable to previous studies that tested on the original SUPERB. However, they indicate current performance levels on these tasks.

No single model excels across all tasks. In phoneme recognition (PR), the SALMONN models were the only ones to achieve a relatively lower phoneme error rate (PER) compared to the others. For keyword spotting (KS), Qwen-Audio-Chat was the only model to perform slightly better than a random guess (50%). Whisper-LLaMA surpassed all other models in intent classification (IC). Regarding emotion recognition (ER), Qwen-Audio-Chat and Qwen2-Audio-7B-Instruct stood out, with the latter even achieving about 75% accuracy. In ASR, SALMONN-13B and WavLLM are the only two models that achieved a word error rate (WER) lower than 10%. Query-by-example (QbE) appears to be very challenging for all models, as none performed better than a random guess (50%). Slot-filling (SF) is evaluated using two metrics: F1 for slot type and CER for slot value. While Whisper-LLaMA and SALMONN-7B are the top two in this task, their results are still far from ideal. In speaker verification (SV), WavLLM was the only model that performed slightly better than random guessing (50%). In speaker diarization (SD), all models performed poorly, with some even reaching a 100% N/A rate. Comparing Figure 3 and Table 2, we observed that performance on core tasks sometimes deviates from trends observed at the domain level. For example, although WavLLM outperforms other models in SV, it lags behind Qwen2-Audio-7B-Chat in speaker and language tasks, revealing its limited capabilities and highlighting the need to enhance its generalizability. In summary, among tested models, some excel in specific tasks, while others perform poorly across the board, and none dominate across all tasks. We believe these findings provide valuable insights for researchers aiming to improve models, starting from core tasks and progressively tackling each domain in the benchmark for broader applicability.

## 6 CONCLUSIONS

This paper presents Dynamic-SUPERB *Phase-2*, the largest benchmark for evaluating instruction-based universal spoken language models, building upon collaborative efforts across the research community. Dynamic-SUPERB covers a wide range of diverse tasks and offers a fine-grained task taxonomy. The recent models show good performance on specific tasks but poor generalization across common tasks like those in SUPERB, highlighting the need for further research on universal models. All materials are made openly available to facilitate reproduction and benchmarking. We sincerely invite researchers to join our vibrant community and collaborate to advance the field.

**Limitations:** Although Dynamic-SUPERB *Phase-2* is the largest and most comprehensive benchmark, we acknowledge its limitations. It lacks comprehensive speech-generation tasks, as *Phase-2* focused on understanding tasks due to the few universal generation models. Despite our efforts to develop the task taxonomy scientifically, new domains may emerge as the benchmark grows, and tasks can be categorized in various ways. While our automatic evaluation pipeline using LLMs correlates well with human evaluations (Appendix E) for current tasks, it may not generalize to all future tasks. Upholding the core spirit of the Dynamic-SUPERB project, we are striving to address these issues and enhance the benchmark for the next phase.

## 7 ETHICS STATEMENT

Dynamic-SUPERB *Phase-2* includes several datasets, and we asked contributors to describe how they used them. Most contributors derived task data by uniformly sampling from the original datasets without applying any special processing. Consequently, Dynamic-SUPERB *Phase-2* inevitably inherits the biases present in these datasets across tasks. Additionally, we asked all contributors to provide the necessary license information to ensure that we can include them in the benchmark properly. More specifically, all task data in Dynamic-SUPERB are derived from datasets with licenses that permit remixing and redistribution. As we expect the benchmark to grow dynamically in the future, we maintain a complete list of dataset licenses on our official GitHub page ([https://github.com/dynamic-superb/dynamic-superb/blob/main/docs/dataset\\_license.md](https://github.com/dynamic-superb/dynamic-superb/blob/main/docs/dataset_license.md)), which will be updated as new tasks are introduced. If there are any changes to the dataset license in the future, we may adjust the benchmark task accordingly.

## 8 REPRODUCIBILITY STATEMENT

We open-source all task data and the evaluation pipeline at <https://github.com/dynamic-superb/dynamic-superb>. For LLM-based evaluation, we set GPT-4o’s temperature to 0 to ensure a stable evaluation process. However, we acknowledge that GPT-4o’s behavior may change with future updates. Therefore, we provide a comprehensive set of evaluation results using LLaMA-3.1-8b-Instruct in Appendix C (Table 6). Additionally, in Appendix E, we present a preliminary study on using open-source LLMs for evaluation, demonstrating promising correlations with human annotations. We encourage researchers to use these open-source LLMs when evaluating their own models to support reproducibility.

## REFERENCES

- Librispeech word. [https://huggingface.co/datasets/speech31/Librispeech\\_word](https://huggingface.co/datasets/speech31/Librispeech_word). Accessed: 2024-10-02.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Eneko Agirre, Daniel Cer, Mona Diab, and Aitor Gonzalez-Agirre. SemEval-2012 task 6: A pilot on semantic textual similarity. In Eneko Agirre, Johan Bos, Mona Diab, Suresh Manandhar, Yuval Marton, and Deniz Yuret (eds.), *\*SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pp. 385–393, Montréal, Canada, 7-8 June 2012. Association for Computational Linguistics. URL <https://aclanthology.org/S12-1051>.
- Andrea Agostinelli, Timo I Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, et al. Musiclm: Generating music from text. *arXiv preprint arXiv:2301.11325*, 2023.
- Afroz Ahamad, Ankit Anand, and Pranesh Bhargava. Accentdb: A database of non-native english accents to assist neural speech recognition. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pp. 5353–5360, Marseille, France, May 2020. European Language Resources Association. URL <https://www.aclweb.org/anthology/2020.lrec-1.659>.
- Anna Aljanaki, Frans Wiering, and Remco C. Veltkamp. Studying emotion induced by music through a crowdsourcing game. *Information Processing & Management*, 52(1):115–128, 2016. ISSN 0306-4573. doi: <https://doi.org/10.1016/j.ipm.2015.03.004>. URL <https://www.sciencedirect.com/science/article/pii/S0306457315000424>. Emotion and Sentiment in Social and Expressive Media.

- Akshay Anantapadmanabhan, Ashwin Bellur, and Hema A Murthy. Modal analysis and transcription of strokes of the mridangam using non-negative matrix factorization. In *2013 IEEE international conference on acoustics, speech and signal processing*, pp. 181–185. IEEE, 2013.
- Xavier Anguera, Luis-J. Rodriguez-Fuentes, Andi Buzo, Florian Metze, Igor Szöke, and Mikel Penagarikano. Quesst2014: Evaluating query-by-example speech search in a zero-resource setting with real-life queries. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5833–5837, 2015. doi: 10.1109/ICASSP.2015.7179090.
- Anthropic. Claude. <https://www.anthropic.com>, 2023. Large language model.
- Rosana Ardila, Megan Branson, Kelly Davis, Michael Kohler, Josh Meyer, Michael Henretty, Reuben Moraes, Lindsay Saunders, Francis Tyers, and Gregor Weber. Common voice: A massively-multilingual speech corpus. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pp. 4218–4222, 2020.
- Muhammad Asif, Muhammad Usaid, Munaf Rashid, Tabarka Rajab, Samreen Hussain, and Sarwar Wasi. Large-scale audio dataset for emergency vehicle sirens and road noises. *Scientific data*, 9(1):599, 2022.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Jon Barker, Shinji Watanabe, Emmanuel Vincent, and Jan Trmal. The Fifth ‘CHiME’ Speech Separation and Recognition Challenge: Dataset, Task and Baselines. In *Proc. Interspeech 2018*, pp. 1561–1565, 2018. doi: 10.21437/Interspeech.2018-1768.
- Emanuele Bastianelli, Andrea Vanzo, Pawel Swietojanski, and Verena Rieser. SLURP: A Spoken Language Understanding Resource Package. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020.
- Sören Becker, Johanna Vielhaben, Marcel Ackermann, Klaus-Robert Müller, Sebastian Lapuschkin, and Wojciech Samek. Audiomnist: Exploring explainable artificial intelligence for audio analysis on a simple benchmark. *Journal of the Franklin Institute*, 361(1):418–428, 2024.
- Jordan J Bird and Ahmad Lotfi. Real-time detection of ai-generated speech for deepfake voice conversion. *arXiv preprint arXiv:2308.12734*, 2023.
- Yonatan Bitton, Hritik Bansal, Jack Hessel, Rulin Shao, Wanrong Zhu, Anas Awadalla, Josh Gardner, Rohan Taori, and Ludwig Schimdt. Visit-bench: a benchmark for vision-language instruction following inspired by real-world use. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, pp. 26898–26922, 2023.
- Dmitry Bogdanov, Minz Won, Philip Tovstogan, Alastair Porter, and Xavier Serra. The mtg-jamendo dataset for automatic music tagging. *ICML*, 2019.
- P. Cano, N. Wack, and P. Herrera. ISMIR04 Genre Identification task dataset (1.0), 2018. URL <https://doi.org/10.5281/zenodo.1302992>.
- Houwei Cao, David G. Cooper, Michael K. Keutmann, Ruben C. Gur, Ani Nenkova, and Ragini Verma. Crema-d: Crowd-sourced emotional multimodal actors dataset. *IEEE Transactions on Affective Computing*, 5(4):377–390, 2014. doi: 10.1109/TAFFC.2014.2336244.
- Jean Carletta, Simone Ashby, Sebastien Bourban, Mike Flynn, Mael Guillemot, Thomas Hain, Jaroslav Kadlec, Vasilis Karaiskos, Wessel Kraaij, Melissa Kronenthal, et al. The ami meeting corpus: A pre-announcement. In *International workshop on machine learning for multimodal interaction*, pp. 28–39. Springer, 2005.
- Santiago Castro, Devamanyu Hazarika, Verónica Pérez-Rosas, Roger Zimmermann, Rada Mihalcea, and Soujanya Poria. Towards multimodal sarcasm detection (an ‘Obviously’ perfect paper). In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4619–4629, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1455. URL <https://aclanthology.org/P19-1455>.

- CEEC. High school english listening exam. <https://www.ceec.edu.tw/>. Accessed: 2024-11-08.
- Liang-Yu CHEN and Jyh-Shing Roger JANG. Stress detection of english words for a capt system using word-length dependent gmm-based bayesian classifiers. *Interdisciplinary Information Sciences*, 18(2):65–70, 2012. doi: 10.4036/iis.2012.65.
- Cheng-Han Chiang and Hung-Yi Lee. Can large language models be an alternative to human evaluations? In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15607–15631, 2023.
- Eleanor Chodroff, Blaž Pažon, Annie Baker, and Steven Moran. Phonetic segmentation of the UCLA phonetics lab archive. In Nicoletta Calzolari, Min-Yen Kan, Veronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue (eds.), *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pp. 12724–12733, Torino, Italia, May 2024. ELRA and ICCL. URL <https://aclanthology.org/2024.lrec-main.1114>.
- Soonbeom Choi, Won Il Kim, Sae Byul Park, Sangeon Yong, and Juhan Nam. Children’s song dataset for singing voice research. In *The 21th International Society for Music Information Retrieval Conference (ISMIR)*. International Society for Music Information Retrieval, 2020.
- Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. *arXiv preprint arXiv:2311.07919*, 2023.
- Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, et al. Qwen2-audio technical report. *arXiv preprint arXiv:2407.10759*, 2024.
- Joris Cosentino, Manuel Pariente, Samuele Cornell, Antoine Deleforge, and Emmanuel Vincent. Librimix: An open-source dataset for generalizable speech separation, 2020.
- Maureen de Seyssel, Marvin Lavechin, Hadrien Titeux, Arthur Thomas, Gwendal Virlet, Andrea Santos Revilla, Guillaume Wisniewski, Bogdan Ludusan, and Emmanuel Dupoux. Prosaudit, a prosodic benchmark for self-supervised speech models. In *INTERSPEECH 2023*, pp. 2963–2967, 2023. doi: 10.21437/Interspeech.2023-438.
- DeepContractor. Musical instrument chord classification, 2024. URL <https://www.kaggle.com/datasets/deepcontractor/musical-instrument-chord-classification>. Accessed: 2024-10-01.
- Michaël Defferrard, Kirell Benzi, Pierre Vandergheynst, and Xavier Bresson. Fma: A dataset for music analysis. In *18th International Society for Music Information Retrieval Conference*, 2017.
- Héctor Delgado, Massimiliano Todisco, Md Sahidullah, Nicholas Evans, Tomi Kinnunen, Kong Aik Lee, and Junichi Yamagishi. Asvspoof 2017 version 2.0: meta-data analysis and baseline enhancements. In *The Speaker and Language Recognition Workshop (Odyssey 2018)*, pp. 296–303, 2018. doi: 10.21437/Odyssey.2018-42.
- Soham Deshmukh, Shuo Han, Hazim Bukhari, Benjamin Elizalde, Hannes Gamper, Rita Singh, and Bhiksha Raj. Audio entailment: Assessing deductive reasoning for audio understanding. *arXiv preprint arXiv:2407.18062*, 2024.
- Konstantinos Drossos, Samuel Lipping, and Tuomas Virtanen. Clotho: an audio captioning dataset. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 736–740, 2020. doi: 10.1109/ICASSP40776.2020.9052990.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

- Ewan Dunbar, Mathieu Bernard, Nicolas Hamilakis, Tu Anh Nguyen, Maureen de Seyssel, Patricia Rozé, Morgane Rivière, Eugene Kharitonov, and Emmanuel Dupoux. The Zero Resource Speech Challenge 2021: Spoken Language Modelling. In *INTERSPEECH 2021*. doi: 10.21437/Interspeech.2021-1755.
- Ewan Dunbar, Nicolas Hamilakis, and Emmanuel Dupoux. Self-supervised language learning from raw audio: Lessons from the zero resource speech challenge. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1211–1226, 2022.
- Jesse Engel, Cinjon Resnick, Adam Roberts, Sander Dieleman, Mohammad Norouzi, Douglas Eck, and Karen Simonyan. Neural audio synthesis of musical notes with wavenet autoencoders. In *International Conference on Machine Learning*, pp. 1068–1077. PMLR, 2017.
- Solène Evain, Ha Nguyen, Hang Le, Marcely Zanon Boito, Salima Mdhaffar, Sina Alisamir, Ziyi Tong, Natalia Tomashenko, Marco Dinarelli, Titouan Parcollet, Alexandre Allauzen, Yannick Estève, Benjamin Lecouteux, François Portet, Solange Rossato, Fabien Ringeval, Didier Schwab, and Laurent Besacier. Lebenchmark: A reproducible framework for assessing self-supervised representation learning from speech. In *Interspeech 2021*, pp. 1439–1443, 2021. doi: 10.21437/Interspeech.2021-556.
- Francesco Foscari, Andrew McLeod, Philippe Rigaux, Florent Jacquemard, and Masahiko Sakai. ASAP: a dataset of aligned scores and performances for piano transcription. In *International Society for Music Information Retrieval Conference (ISMIR)*, pp. 534–541, 2020.
- Joel Frank and Lea Schönherr. Wavefake: A data set to facilitate audio deepfake detection. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Daniel Galvez, Greg Diamos, Juan Torres, Keith Achorn, Juan Cerón, Anjali Gopi, David Kanter, Max Lam, Mark Mazumder, and Vijay Janapa Reddi. The people’s speech: A large-scale diverse english speech recognition dataset for commercial usage. In J. Vanschoren and S. Yeung (eds.), *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1, 2021. URL [https://datasets-benchmarks-proceedings.neurips.cc/paper\\_files/paper/2021/file/202cb962ac59075b964b07152d234b70-Paper-round1.pdf](https://datasets-benchmarks-proceedings.neurips.cc/paper_files/paper/2021/file/202cb962ac59075b964b07152d234b70-Paper-round1.pdf).
- Mandar Gogate, Kia Dashtipour, Ahsan Adeel, and Amir Hussain. Cochleanet: A robust language-independent audio-visual model for real-time speech enhancement. *Information Fusion*, 63:273–285, 2020.
- Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *circulation*, 101(23):e215–e220, 2000.
- Yuan Gong, Jin Yu, and James Glass. Vocalsound: A dataset for improving human vocal sounds recognition. In *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 151–155, 2022. doi: 10.1109/ICASSP43922.2022.9746828.
- Yuan Gong, Alexander H Liu, Hongyin Luo, Leonid Karlinsky, and James Glass. Joint audio and speech understanding. In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pp. 1–8. IEEE, 2023.
- Zhifang Guo, Yichong Leng, Yihan Wu, Sheng Zhao, and Xu Tan. Promptts: Controllable text-to-speech with text descriptions. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.
- Curtis Hawthorne, Andriy Stasyuk, Adam Roberts, Ian Simon, Cheng-Zhi Anna Huang, Sander Dieleman, Erich Elsen, Jesse Engel, and Douglas Eck. Enabling factorized piano music modeling and generation with the MAESTRO dataset. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=r1lYRjC9F7>.

- Addison Howard et al. Cornell birdcall identification, 2020. URL <https://kaggle.com/competitions/birdsong-recognition>.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021.
- Shujie Hu, Long Zhou, Shujie Liu, Sanyuan Chen, Hongkun Hao, Jing Pan, Xunying Liu, Jinyu Li, Sunit Sivasankaran, Linqun Liu, et al. Wavllm: Towards robust and adaptive speech large language model. *arXiv preprint arXiv:2404.00656*, 2024.
- Chien-yu Huang, Ke-Han Lu, Shih-Heng Wang, Chi-Yuan Hsiao, Chun-Yi Kuan, Haibin Wu, Siddhant Arora, Kai-Wei Chang, Jiatong Shi, Yifan Peng, et al. Dynamic-superb: Towards a dynamic, collaborative, and comprehensive instruction-tuning benchmark for speech. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 12136–12140. IEEE, 2024a.
- Kuan-Po Huang, Chih-Kai Yang, Yu-Kuan Fu, Ewan Dunbar, and Hung-yi Lee. Zero resource code-switched speech benchmark using speech utterance pairs for multiple spoken languages. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 10006–10010. IEEE, 2024b.
- Eric Humphrey, Simon Durand, and Brian McFee. Openmic-2018: An open data-set for multiple instrument recognition. In *ISMIR*, pp. 438–444, 2018.
- IEEE DCASE 2016 Challenge. IEEE DCASE 2016 Challenge. <http://www.cs.tut.fi/sgn/arg/dcase2016/>, 2016. Accessed: 2025-03-04.
- Keith Ito and Linda Johnson. The lj speech dataset. <https://keithito.com/LJ-Speech-Dataset/>, 2017.
- Tahir Javed, Kaushal Bhogale, Abhigyan Raman, Pratyush Kumar, Anoop Kunchukuttan, and Mitesh M Khapra. Indicsuperb: A speech processing universal performance benchmark for indian languages. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 12942–12950, 2023.
- Denys Katerenchuk, David Guy Brizan, and Andrew Rosenberg. Interpersonal relationship labels for the CALLHOME corpus. In Nicoletta Calzolari, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Koiti Hasida, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Asuncion Moreno, Jan Odijk, Stelios Piperidis, and Takenobu Tokunaga (eds.), *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan, May 2018. European Language Resources Association (ELRA). URL <https://aclanthology.org/L18-1592>.
- Tyler Kendall and Charlie Farrington. The corpus of regional african american language, 2023. URL <https://doi.org/10.7264/1ad5-6t35>.
- Bongjun Kim, Madhav Ghei, Bryan Pardo, and Zhiyao Duan. Vocal imitation set: a dataset of vocally imitated sound events using the audioset ontology. In *DCASE*, pp. 148–152, 2018.
- Peter Knees, Ángel Faraldo Pérez, Herrera Boyer, Richard Vogl, Sebastian Böck, Florian Hörschläger, Mickael Le Goff, et al. Two data sets for tempo estimation and key detection in electronic dance music annotated from user corrections. In *Proceedings of the 16th International Society for Music Information Retrieval Conference (ISMIR); 2015 Oct 26-30; Málaga, Spain.[Málaga]: International Society for Music Information Retrieval, 2015. p. 364-70. International Society for Music Information Retrieval (ISMIR), 2015*.
- Tom Ko, Vijayaditya Peddinti, Daniel Povey, Michael L. Seltzer, and Sanjeev Khudanpur. A study on data augmentation of reverberant speech for robust speech recognition. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5220–5224, 2017. doi: 10.1109/ICASSP.2017.7953152.

- Edith Law, Kris West, Michael I Mandel, Mert Bay, and J Stephen Downie. Evaluation of algorithms using games: The case of music tagging. In *ISMIR*, pp. 387–392. Citeseer, 2009.
- Colin Lea, Vikramjit Mitra, Aparna Joshi, Sachin Kajarekar, and Jeffrey P. Bigham. Sep-28k: A dataset for stuttering event detection from podcasts with people who stutter. In *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6798–6802, 2021. doi: 10.1109/ICASSP39728.2021.9413520.
- Keon Lee, Kyumin Park, and Daeyoung Kim. Dailytalk: Spoken dialogue dataset for conversational text-to-speech. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pp. 19730–19742. PMLR, 2023a.
- Yizhi LI, Ruibin Yuan, Ge Zhang, Yinghao Ma, Xingran Chen, Hanzhi Yin, Chenghao Xiao, Chenghua Lin, Anton Ragni, Emmanouil Benetos, Norbert Gyenge, Roger Dannenberg, Ruibo Liu, Wenhua Chen, Gus Xia, Yemin Shi, Wenhao Huang, Zili Wang, Yike Guo, and Jie Fu. MERT: Acoustic music understanding model with large-scale self-supervised training. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=w3YZ9MSlBu>.
- Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*, 2023b.
- Guan-Ting Lin, Prashanth Gurunath Shivakumar, Ankur Gandhe, Chao-Han Huck Yang, Yile Gu, Shalini Ghosh, Andreas Stolcke, Hung-yi Lee, and Ivan Bulyko. Paralinguistics-enhanced large language modeling of spoken dialogue. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 10316–10320. IEEE, 2024.
- Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. Audioldm: text-to-audio generation with latent diffusion models. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 21450–21474, 2023a.
- Shansong Liu, Atin Sakkeer Hussain, Chenshuo Sun, and Ying Shan. Music understanding llama: Advancing text-to-music generation with question answering and captioning. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 286–290. IEEE, 2024.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval: Nlg evaluation using gpt-4 with better human alignment. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 2511–2522, 2023b.
- Steven R Livingstone and Frank A Russo. The ryerson audio-visual database of emotional speech and song (ravdess): A dynamic, multimodal set of facial and vocal expressions in north american english. *PLoS one*, 13(5):e0196391, 2018.
- Holy Lovenia, Samuel Cahyawijaya, Genta Indra Winata, Peng Xu, Xu Yan, Zihan Liu, Rita Frieske, Tiezheng Yu, Wenliang Dai, Elham J Barezi, et al. Ascend: A spontaneous chinese-english dataset for code-switching in multi-turn conversation. In *Proceedings of the 13th Language Resources and Evaluation Conference (LREC)*, 2022.
- Ken MacLean. Voxforge, 2018. Ken MacLean.[Online]. Available: <http://www.voxforge.org/home>. [Acedido em 2012].
- Gallil Maimon, Amit Roth, and Yossi Adi. A suite for acoustic language model evaluation. *arXiv preprint arXiv:2409.07437*, 2024.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. Cross-task generalization via natural language crowdsourcing instructions. In *ACL*, 2022.



- National Chengchi University. The nccu (national chengchi university) corpus of spoken taiwan mandarin, n.d. URL <http://spokentaiwanmandarin.nccu.edu.tw/>. [Accessed: 2024-10-01].
- Inês Nolasco, Alessandro Terenzi, Stefania Cecchi, Simone Orcioni, Helen L Bear, and Emmanouil Benetos. Audio-based identification of beehive states. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 8256–8260. IEEE, 2019.
- Lara Orlandic, Tomas Teijeiro, and David Atienza. The coughvid crowdsourcing dataset, a corpus for the study of large-scale cough analysis algorithms. *Scientific Data*, 8(1):156, 2021.
- Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: an asr corpus based on public domain audio books. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pp. 5206–5210. IEEE, 2015.
- Yagya Raj Pandeya and Joonwhoan Lee. Domestic cat sound classification using transfer learning. *International Journal of Fuzzy Logic and Intelligent Systems*, 18(2):154–160, 2018.
- Yagya Raj Pandeya, Dongwhoon Kim, and Joonwhoan Lee. Domestic cat sound classification using learned features from deep neural nets. *Applied Sciences*, 8(10):1949, 2018.
- Titouan Parcollet, Ha Nguyen, Solene Evain, Marcely Zanon Boito, Adrien Pupier, Salima Mdhafar, Hang Le, Sina Alisamir, Natalia Tomashenko, Marco Dinarelli, et al. Lebenchmark 2.0: a standardized, replicable and enhanced framework for self-supervised representations of french speech. *arXiv preprint arXiv:2309.05472*, 2023.
- Karol J Piczak. Esc: Dataset for environmental sound classification. In *Proceedings of the 23rd ACM international conference on Multimedia*, pp. 1015–1018, 2015.
- Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. MELD: A multimodal multi-party dataset for emotion recognition in conversations. In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 527–536, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1050. URL <https://aclanthology.org/P19-1050>.
- Sathvik Prasad and Bradley Reaves. Robocall audio from the ftc’s project point of no entry.
- Vineel Pratap, Qiantong Xu, Anuroop Sriram, Gabriel Synnaeve, and Ronan Collobert. Mls: A large-scale multilingual dataset for speech research. In *Interspeech 2020*, pp. 2757–2761, 2020. doi: 10.21437/Interspeech.2020-2826.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pp. 28492–28518. PMLR, 2023.
- Zafar Rafii, Antoine Liutkus, Fabian-Robert Stöter, Stylianos Ioannis Mimilakis, and Rachel Bittner. MUSDB18-HQ - an uncompressed version of musdb18, December 2019. URL <https://doi.org/10.5281/zenodo.3338373>.
- Muhammad Mamunur Rashid, Guiqing Li, and Chengrui Du. Nonspeech7k dataset: Classification and analysis of human non-speech sound. *IET Signal Processing*, 17(6):e12233, 2023.
- Matthew Rice, Christian J. Steinmetz, George Fazekas, and Joshua D. Reiss. General purpose audio effect removal. In *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, 2023.
- J. Salamon, C. Jacoby, and J. P. Bello. A dataset and taxonomy for urban sound research. In *22nd ACM International Conference on Multimedia (ACM-MM’14)*, pp. 1041–1044, Orlando, FL, USA, Nov. 2014.
- Miguel Sarabia, Elena Menyaylenko, Alessandro Toso, Skyler Seto, Zakaria Aldeneh, Shadi Pirhoseinloo, Luca Zappella, Barry-John Theobald, Nicholas Apostoloff, and Jonathan Sheaffer. Spatial librispeech: An augmented dataset for spatial audio learning. In *INTERSPEECH 2023*, pp. 3724–3728, 2023. doi: 10.21437/Interspeech.2023-2117.

- Jiatong Shi, Dan Berrebbi, William Chen, En-Pei Hu, Wei-Ping Huang, Ho-Lam Chung, Xuankai Chang, Shang-Wen Li, Abdelrahman Mohamed, Hung yi Lee, and Shinji Watanabe. MI-superb: Multilingual speech universal performance benchmark. In *INTERSPEECH 2023*, pp. 884–888, 2023. doi: 10.21437/Interspeech.2023-1316.
- Jiatong Shi, Yueqian Lin, Xinyi Bai, Keyi Zhang, Yuning Wu, Yuxun Tang, Yifeng Yu, Qin Jin, and Shinji Watanabe. Singing voice data scaling-up: An introduction to ace-opencpop and ace-kising. *arXiv preprint arXiv:2401.17619*, 2024a.
- Jiatong Shi, Shih-Heng Wang, William Chen, Martijn Bartelds, Vanya Bannihatti Kumar, Jinchuan Tian, Xuankai Chang, Dan Jurafsky, Karen Livescu, Hung yi Lee, and Shinji Watanabe. MI-superb 2.0: Benchmarking multilingual speech models across modeling constraints, languages, and datasets. In *Interspeech 2024*, pp. 1230–1234, 2024b. doi: 10.21437/Interspeech.2024-2248.
- Kazuki Shimada, Archontis Politis, Parthasaarathy Sudarsanam, Daniel A Krause, Kengo Uchida, Sharath Adavanne, Aapo Hakala, Yuichiro Koyama, Naoya Takahashi, Shusuke Takahashi, et al. Starss23: An audio-visual dataset of spatial recordings of real scenes with spatiotemporal annotations of sound events. *Advances in Neural Information Processing Systems*, 36, 2024.
- Suwon Shon, Ankita Pasad, Felix Wu, Pablo Brusco, Yoav Artzi, Karen Livescu, and Kyu J Han. Slue: New benchmark tasks for spoken language understanding evaluation on natural speech. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7927–7931. IEEE, 2022.
- Suwon Shon, Siddhant Arora, Chyi-Jiunn Lin, Ankita Pasad, Felix Wu, Roshan S Sharma, Wei-Lun Wu, Hung-Yi Lee, Karen Livescu, and Shinji Watanabe. Slue phase-2: A benchmark suite of diverse spoken language understanding tasks. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8906–8937, 2023.
- David Snyder, Guoguo Chen, and Daniel Povey. Musan: A music, speech, and noise corpus. *arXiv preprint arXiv:1510.08484*, 2015.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*, 2022.
- Fabian-Robert Stöter, Soumitro Chakrabarty, Emanuël Habets, and Bernd Edler. LibriCount, a dataset for speaker count estimation, 2018. URL <https://doi.org/10.5281/zenodo.1216072>.
- Dan Stowell, Michael D Wood, Hanna Pamuła, Yannis Stylianou, and Hervé Glotin. Automatic acoustic detection of birds through deep learning: the first bird audio detection challenge. *Methods in Ecology and Evolution*, 10(3):368–380, 2019.
- Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, MA Zejun, and Chao Zhang. Salmonn: Towards generic hearing abilities for large language models. In *The Twelfth International Conference on Learning Representations*, 2024a.
- Yuxun Tang, Jiatong Shi, Yuning Wu, and Qin Jin. Singmos: An extensive open-source singing voice dataset for mos prediction. *arXiv preprint arXiv:2406.10911*, 2024b.
- Mi Tian, Ajay Srinivasamurthy, Mark Sandler, and Xavier Serra. A study of instrument-wise onset detection in beijing opera percussion ensembles. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2159–2163. IEEE, 2014.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.

- Sheng-Chi Tsai. English lexical stress detection and sentence-based intonation assessment based on contour shape description. Master's thesis, National Taiwan University, Taipei, Taiwan, 2015. Available at <https://hdl.handle.net/11296/jjzzak>.
- Ching-Yu Tseng. An initial study on stress detection for spoken english. Master's thesis, National Tsing Hua University, Hsinchu, Taiwan, 2008. Available at <https://hdl.handle.net/11296/qr88s2>.
- Joseph Turian, Jordie Shier, Humair Raj Khan, Bhiksha Raj, Björn W Schuller, Christian J Steinmetz, Colin Malloy, George Tzanetakis, Gissel Velarde, Kirk McNally, et al. Hear: Holistic evaluation of audio representations. In *NeurIPS 2021 Competitions and Demonstrations Track*, pp. 125–145. PMLR, 2022.
- Nicolas Turpault, Romain Serizel, Ankit Parag Shah, and Justin Salamon. Sound event detection in domestic environments with weakly labeled data and soundscape synthesis. In *Workshop on Detection and Classification of Acoustic Scenes and Events*, 2019.
- Jörgen Valk and Tanel Alumäe. Voxlingua107: a dataset for spoken language recognition. In *2021 IEEE Spoken Language Technology Workshop (SLT)*, pp. 652–658. IEEE, 2021.
- Arun Balajee Vasudevan, Dengxin Dai, and Luc Van Gool. Semantic object prediction and spatial sound super-resolution with binaural sounds. In *European Conference on Computer Vision*, pp. 638–655. Springer, 2020.
- Changhan Wang, Anne Wu, Jiatao Gu, and Juan Pino. Covost 2 and massively multilingual speech translation. In *Interspeech 2021*, pp. 2247–2251, 2021. doi: 10.21437/Interspeech.2021-2027.
- Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. Is chatgpt a good nlg evaluator? a preliminary study. In *Proceedings of EMNLP Workshop*, pp. 1, 2023.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al. Super-naturalinstructions: generalization via declarative instructions on 1600+ tasks. In *EMNLP*, 2022.
- Pete Warden. Speech commands: A dataset for limited-vocabulary speech recognition. *arXiv preprint arXiv:1804.03209*, 2018.
- Shinji Watanabe, Michael Mandel, Jon Barker, Emmanuel Vincent, Ashish Arora, Xuankai Chang, Sanjeev Khudanpur, Vimal Manohar, Daniel Povey, Desh Raj, et al. Chime-6 challenge: Tackling multispeaker speech recognition for unsegmented recordings. In *CHiME 2020-6th International Workshop on Speech Processing in Everyday Environments*, 2020.
- WavSource. Sound effects library. <https://www.wavsource.com>, n.d. Accessed: 2024-10-01.
- whats2000. Human screaming detection dataset, 2023. URL <https://www.kaggle.com/datasets/whats2000/human-screaming-detection-dataset/data>. Accessed: 2024-10-02.
- Gordon Wichern, Joe Antognini, Michael Flynn, Licheng Richard Zhu, Emmett McQuinn, Dwight Crow, Ethan Manilow, and Jonathan Le Roux. Wham!: Extending speech separation to noisy environments. In *Interspeech 2019*, pp. 1368–1372, 2019. doi: 10.21437/Interspeech.2019-2821.
- Julia Wilkins, Prem Seetharaman, Alison Wahl, and Bryan Pardo. Vocalset: A singing voice dataset. In *ISMIR*, pp. 468–474, 2018.
- Zhizheng Wu, Tomi Kinnunen, Nicholas Evans, Junichi Yamagishi, Cemal Haniç, Md. Sahidullah, and Aleksandr Sizov. Asvspoof 2015: the first automatic speaker verification spoofing and countermeasures challenge. In *Interspeech 2015*, pp. 2037–2041, 2015. doi: 10.21437/Interspeech.2015-462.

- Junichi Yamagishi, Christophe Veaux, and Kirsten MacDonald. CSTR VCTK Corpus: English Multi-speaker Corpus for CSTR Voice Cloning Toolkit (version 0.92). sound, 2019. URL <https://doi.org/10.7488/ds/2645>.
- Chih-Kai Yang, Kuan-Po Huang, Ke-Han Lu, Chun-Yi Kuan, Chi-Yuan Hsiao, and Hung-yi Lee. Investigating zero-shot generalizability on mandarin-english code-switched asr and speech-to-text translation of recent foundation models with self-supervision and weak supervision. In *2024 IEEE International Conference on Acoustics, Speech, and Signal Processing Workshops (ICASSPW)*, pp. 540–544. IEEE, 2024a.
- Qian Yang, Jin Xu, Wenrui Liu, Yunfei Chu, Ziyue Jiang, Xiaohuan Zhou, Yichong Leng, Yuanjun Lv, Zhou Zhao, Chang Zhou, et al. Air-bench: Benchmarking large audio-language models via generative comprehension. *arXiv preprint arXiv:2402.07729*, 2024b.
- 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. Superb: Speech processing universal performance benchmark. In *Interspeech 2021*, pp. 1194–1198, 2021. doi: 10.21437/Interspeech.2021-1775.
- Qinyuan Ye, Bill Yuchen Lin, and Xiang Ren. Crossfit: A few-shot learning challenge for cross-task generalization in nlp. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 7163–7189, 2021.
- Ruibin Yuan, Yinghao Ma, Yizhi Li, Ge Zhang, Xingran Chen, Hanzhi Yin, Yiqi Liu, Jiawen Huang, Zeyue Tian, Binyue Deng, et al. Marble: Music audio representation benchmark for universal evaluation. *Advances in Neural Information Processing Systems*, 36:39626–39647, 2023.
- Yongyi Zang, Jiatong Shi, You Zhang, Ryuichi Yamamoto, Jionghao Han, Yuxun Tang, Shengyuan Xu, Wenxiao Zhao, Jing Guo, Tomoki Toda, and Zhiyao Duan. Ctrsvdd: A benchmark dataset and baseline analysis for controlled singing voice deepfake detection. In *Interspeech 2024*, pp. 4783–4787, 2024. doi: 10.21437/Interspeech.2024-2242.
- Heiga Zen, Viet Dang, Rob Clark, Yu Zhang, Ron J. Weiss, Ye Jia, Zhifeng Chen, and Yonghui Wu. Libritts: A corpus derived from librispeech for text-to-speech. In *Interspeech 2019*, pp. 1526–1530, 2019. doi: 10.21437/Interspeech.2019-2441.
- Junbo Zhang, Zhiwen Zhang, Yongqing Wang, Zhiyong Yan, Qiong Song, Yukai Huang, Ke Li, Daniel Povey, and Yujun Wang. speechocean762: An open-source non-native english speech corpus for pronunciation assessment. In *Interspeech 2021*, pp. 3710–3714, 2021. doi: 10.21437/Interspeech.2021-1259.
- Lichao Zhang, Ruiqi Li, Shoutong Wang, Liquan Deng, Jinglin Liu, Yi Ren, Jinzheng He, Rongjie Huang, Jieming Zhu, Xiao Chen, and Zhou Zhao. M4singer: A multi-style, multi-singer and musical score provided mandarin singing corpus. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. URL <https://openreview.net/forum?id=qiDmAaG6mP>.
- Siqi Zheng, Luyao Cheng, Yafeng Chen, Hui Wang, and Qian Chen. 3d-speaker: A large-scale multi-device, multi-distance, and multi-dialect corpus for speech representation disentanglement. *arXiv preprint arXiv:2306.15354*, 2023.

## A TASK TAXONOMY

Figures 4 and 5 show the complete task taxonomy in Dynamic-SUPERB *Phase-2*. Additionally, Table 3 compares our task taxonomy with the INTERSPEECH conference and EDICS from IEEE SPS. While not identical, our taxonomy is well-aligned with those organized by professional conferences and institutions.

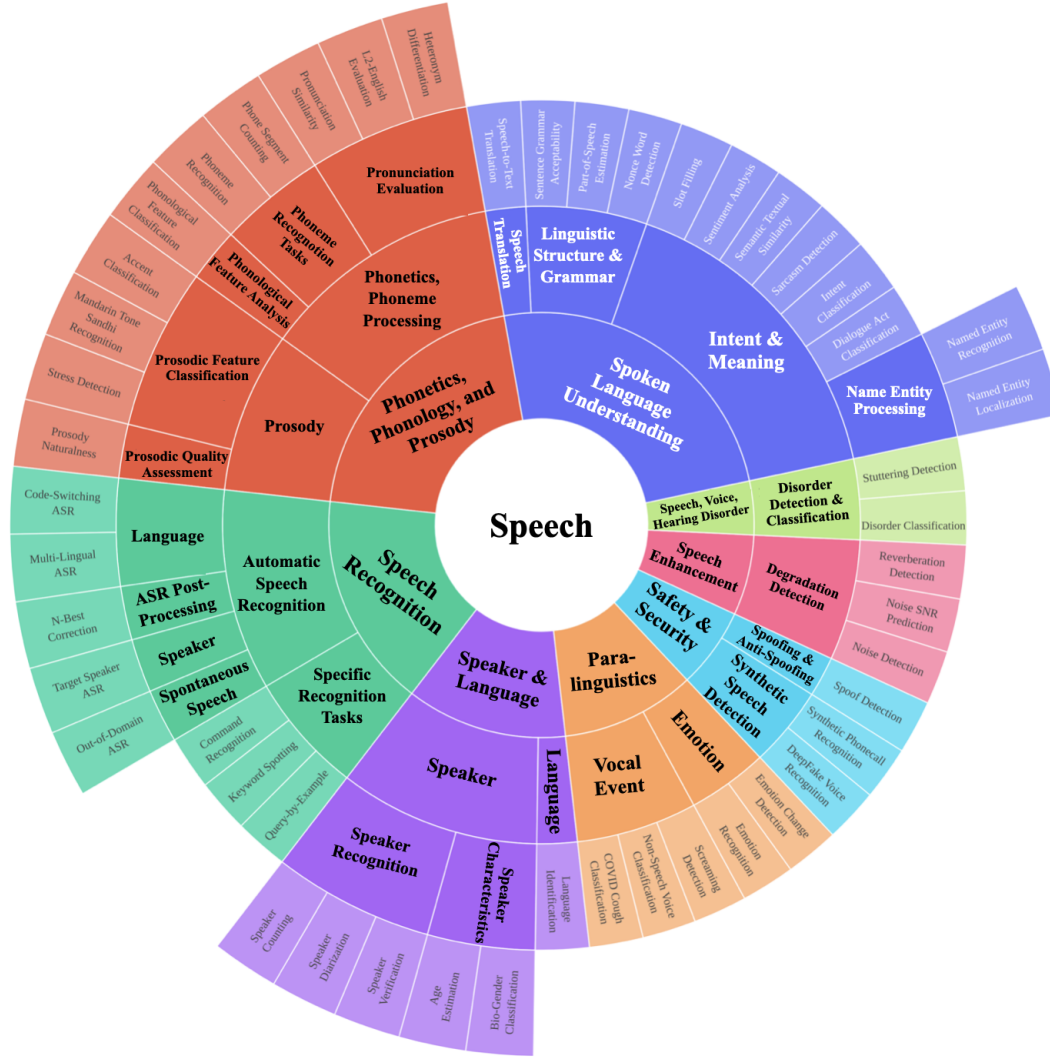


Figure 4: Task taxonomy of speech.

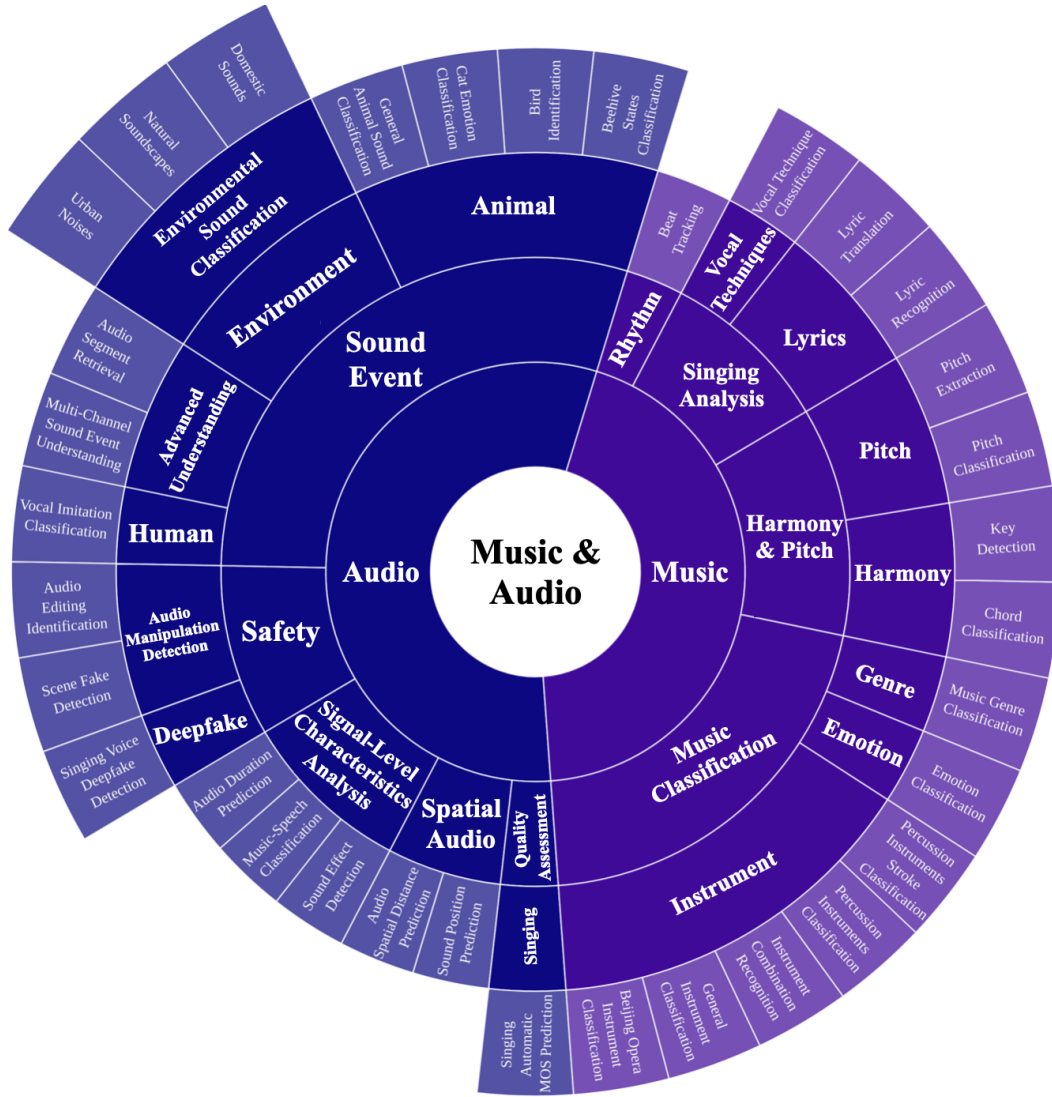


Figure 5: Task taxonomy of audio and music.

Table 3: Comparison of task taxonomy in Dynamic-SUPERB *Phase-2* with INTERSPEECH and EDICS. A “✓” indicates that a conference session, area, or EDICS category matches or closely aligns with the taxonomy.

| <b>Taxonomy</b>                        | <b>INTERSPEECH</b> | <b>EDICS</b> |
|--|--------------------|--------------|
| Music/Harmony & Pitch                  |                    | ✓            |
| Music/Music Classification             |                    |              |
| Music/Rhythm Analysis                  |                    | ✓            |
| Audio/Quality Assessment               |                    | ✓            |
| Audio/Safety                           |                    | ✓            |
| Audio/Signal-Characteristics Analysis  |                    |              |
| Audio/Singing Analysis                 | ✓                  |              |
| Audio/Sound Event                      | ✓                  | ✓            |
| Audio/Spatial Audio Analysis           | ✓                  | ✓            |
| Speech/Paralinguistics                 | ✓                  |              |
| Speech/Phonetics, Phonology, Prosody   | ✓                  |              |
| Speech/Safety and Security             | ✓                  |              |
| Speech/Speaker & Language              | ✓                  |              |
| Speech/Speech Enhancement              | ✓                  | ✓            |
| Speech/Speech Recognition              | ✓                  | ✓            |
| Speech/Speech, Voice, Hearing Disorder | ✓                  |              |
| Speech/Spoken Language Understanding   | ✓                  | ✓            |

## B DYNAMIC-SUPERB TASKS

Table 4: The list of all tasks in Dynamic-SUPERB. Tasks are ordered by the taxonomy.

| <b>Domain</b>                          | <b>Task</b>                                      | <b>Dataset</b>   |
|--|--|--|
| Audio / Harmony & Pitch / Harmony      | Chord Classification                             | Musical Instrument Chord Classification (Audio) (DeepContractor, 2024) |
| Audio / Harmony & Pitch / Harmony      | MARBLE Key Detection                             | GiantSteps (key) (Knees et al., 2015)                                  |
| Audio / Harmony & Pitch / Pitch        | HEAR Music Transcription                         | MAESTRO (Hawthorne et al., 2019)                                       |
| Audio / Harmony & Pitch / Pitch        | HEAR Percussion Instruments Tonic Classification | Mridangam Stroke (Anantapadmanabhan et al., 2013)                      |
| Audio / Harmony & Pitch / Pitch        | Instrument Pitch Classification                  | Nsynth (Engel et al., 2017)  |
| Audio / Harmony & Pitch / Pitch        | Pitch Extraction By Lyrics                       | CSD (Choi et al., 2020)  |
| Audio / Music Classification / Emotion | Emotion Classification In Songs                  | EMOTIFY (Aljanaki et al., 2016)  |
| Audio / Music Classification / Emotion | MARBLE Emotion Detection                         | MTG-Jamendo (Bogdanov et al., 2019)                                    |
| Audio / Music Classification / Genre   | HEAR Music Genre Classification                  | ISMIR04 (Cano et al., 2018)  |

*Continued on next page*

| Domain                                    | Task  | Dataset  |
|---|---|--|
| Audio / Music Classification / Genre      | MARBLE Genre Classification<br>MTG-Genre          | MTG-Jamendo (Bogdanov et al., 2019)  |
| Audio / Music Classification / Genre      | MARBLE Music Tagging                              | MagnaTagATune (Law et al., 2009)   |
| Audio / Music Classification / Genre      | MARBLE Music Tagging                              | MTG-Jamendo (Bogdanov et al., 2019)  |
| Audio / Music Classification / Genre      | Music Genre Classification                        | FMA (Defferrard et al., 2017)  |
| Audio / Music Classification / Instrument | HEAR Percussion Instruments Classification        | Beijing Opera Percussion Instrument Dataset (Tian et al., 2014)                                  |
| Audio / Music Classification / Instrument | HEAR Percussion Instruments Stroke Classification | Mridangam Stroke (Anantapadmanabhan et al., 2013)  |
| Audio / Music Classification / Instrument | Instrument Classification                         | Nsynth (Engel et al., 2017)  |
| Audio / Music Classification / Instrument | Instrument Combination Recognition                | OpenMIC-2018 (Humphrey et al., 2018)   |
| Audio / Music Classification / Instrument | Instrument Source Classification                  | Nsynth (Engel et al., 2017)  |
| Audio / Music Classification / Instrument | MARBLE Instrument Classification                  | MTG-Jamendo (Bogdanov et al., 2019)  |
| Audio / Quality Assessment / Singing      | Singing Automatic MOS Prediction                  | SingMOS (Tang et al., 2024b)   |
| Audio / Rhythm Analysis                   | MARBLE Beat Tracking ASAP                         | ASAP (Foscarin et al., 2020)   |
| Audio / Rhythm Analysis                   | Music Beat Tracking ASAP                          | ASAP (Foscarin et al., 2020)   |
| Audio / Safety / Audio Integrity          | Audio Editing Identification                      | People’s Speech (Galvez et al., 2021)  |
| Audio / Safety / Audio Integrity          | Scene Fake Detection                              | ASPIRE (Gogate et al., 2020)   |
| Audio / Safety / Deepfake                 | Audio Deep Fake Detection                         | LJSpeech (Ito & Johnson, 2017)<br>WaveFake (Frank & Schönherr)<br>MUSDB18HQ (Rafii et al., 2019) |
| Audio / Safety / Deepfake                 | Singing Voice Deepfake Detection                  | CtrSVDD (Zang et al., 2024)<br>ACEKiSing (Shi et al., 2024a)<br>M4Singer (Zhang et al., 2022)    |
| Audio / Signal Characteristics Analysis   | Audio Duration Prediction                         | NTUML2021 (Yang et al., 2024a)   |

*Continued on next page*



| <b>Domain</b>                                | <b>Task</b>                                | <b>Dataset</b>   |
|--|--|--|
| Audio / Signal Characteristics Analysis      | HEAR Music Speech Classification           | MAESTRO (Hawthorne et al., 2019)<br>Librispeech (Panayotov et al., 2015)     |
| Audio / Signal Characteristics Analysis      | Sound Effect Detection                     | RemFx (Rice et al., 2023)  |
| Audio / Signal Characteristics Analysis      | Speech Detection                           | LibriSpeech-TestClean (Panayotov et al., 2015)                               |
| Audio / Signal Characteristics Analysis      | Speech Detection                           | LibriSpeech-TestOther (Panayotov et al., 2015)                               |
| Audio / Signal Characteristics Analysis      | Speech Detection                           | LJSpeech (Ito & Johnson, 2017)   |
| Audio / Singing Analysis / Lyrics            | Children Song Transcript Verification      | CSD (Choi et al., 2020)  |
| Audio / Singing Analysis / Lyrics            | Lyric Translation                          | SingSet (self-recorded)  |
| Audio / Singing Analysis / Lyrics            | Song Lyric Recognition                     | SingSet (self-recorded)  |
| Audio / Singing Analysis / Vocal Techniques  | MARBLE Vocal Technique Detection           | VocalSet (Wilkins et al., 2018)  |
| Audio / Sound Event / Advanced Understanding | Audio Segment Retrieval                    | Clotho (Drossos et al., 2020)  |
| Audio / Sound Event / Advanced Understanding | HEAR Sound Event Detection                 | DCASE2016Task2 (IEEE DCASE 2016 Challenge, 2016)                             |
| Audio / Sound Event / Advanced Understanding | Multi-channel Sound Event Understanding    | STARSS23 (Shimada et al., 2024)  |
| Audio / Sound Event / Animal                 | Animal Classification                      | WaveSource-Test (WavSource, n.d.)  |
| Audio / Sound Event / Animal                 | Bird Sound Detection                       | Warblrb10k (Stowell et al., 2019)  |
| Audio / Sound Event / Animal                 | Cat Emotion Classification                 | CatSoundClassificationDataset-V2 (Pandeya et al., 2018; Pandeya & Lee, 2018) |
| Audio / Sound Event / Animal                 | Cornell Birdcall Identification            | Cornell Birdcall Identification (Howard et al., 2020)                        |
| Audio / Sound Event / Animal                 | Environmental Sound Classification         | ESC50-Animals (Piczak, 2015)   |
| Audio / Sound Event / Animal                 | HEAR Beehive States Classification         | Beehive States (Nolasco et al., 2019)  |
| Audio / Sound Event / Environment            | Domestic Environment Sound Event Detection | DESED-PublicEval (Turpault et al., 2019)                                     |

*Continued on next page*

| Domain   | Task                                      | Dataset   |
|--|---|---|
| Audio / Sound Event / Environment                | Emergency Traffic Detection               | Large-Scale-Audio-dataset (Asif et al., 2022)                               |
| Audio / Sound Event / Environment                | Environmental Sound Classification        | ESC50-ExteriorAndUrbanNoises (Piczak, 2015)                                 |
| Audio / Sound Event / Environment                | Environmental Sound Classification        | ESC50-InteriorAndDomesticSounds (Piczak, 2015)                              |
| Audio / Sound Event / Environment                | Environmental Sound Classification        | ESC50-NaturalSoundscapesAndWaterSounds (Piczak, 2015)                       |
| Audio / Sound Event / Environment                | Environmental Sound Classification        | UrbanSound8K-UrbanNoises (Salamon et al., 2014)                             |
| Audio / Sound Event / Environment                | Environment Recognition                   | ESC50 (Piczak, 2015)  |
| Audio / Sound Event / Environment                | HEAR Environmental Sound Classification   | ESC50 (Piczak, 2015)  |
| Audio / Sound Event / Human                      | HEAR Vocal Imitation Classification       | Vocal Imitations (Kim et al., 2018)   |
| Audio / Spatial Audio                            | Audio Spatial Distance Prediction         | Spatial LibriSpeech (Sarabia et al., 2023)                                  |
| Audio / Spatial Audio                            | How Far Are You                           | 3DSpeaker (Zheng et al., 2023)  |
| Audio / Spatial Audio                            | Sound Position Prediction                 | BinauralSoundPerception (Vasudevan et al., 2020)                            |
| Speech / Paralinguistics / Emotion Analysis      | Dialogue Emotion Classification           | DailyTalk (Lee et al., 2023)  |
| Speech / Paralinguistics / Emotion Analysis      | Emoji Grounded Speech Emotion Recognition | RAVDESS (Livingstone & Russo, 2018)   |
| Speech / Paralinguistics / Emotion Analysis      | Emotion Change Detection                  | RAVDESS (Livingstone & Russo, 2018)   |
| Speech / Paralinguistics / Emotion Analysis      | Emotion Recognition                       | MELD (Poria et al., 2019)   |
| Speech / Paralinguistics / Emotion Analysis      | HEAR Emotion Recognition                  | CREMAD (Cao et al., 2014)   |
| Speech / Paralinguistics / Emotion Analysis      | SuperbER RAVDESS                          | RAVDESS (Livingstone & Russo, 2018)   |
| Speech / Paralinguistics / Vocal Event Detection | Covid19 Cough Audio Classification        | CoughVid (Orlandic et al., 2021)  |
| Speech / Paralinguistics / Vocal Event Detection | Environmental Sound Classification        | ESC50-HumanAndNonSpeechSounds (Piczak, 2015)                                |
| Speech / Paralinguistics / Vocal Event Detection | Human Non-Speech Sound Recognition        | Nonspeech7k-test (Rashid et al., 2023)<br>CommonVoice (Ardila et al., 2020) |
| Speech / Paralinguistics / Vocal Event Detection | Human Screaming Detection                 | Environmentdb (whats2000, 2023)   |

*Continued on next page*

| Domain  | Task                                | Dataset   |
|---|-------------------------------------|---|
| Speech / Paralinguistics / Vocal Event Detection  | Vocal Sound Recognition             | VocalSound (Gong et al., 2022)                                  |
| Speech / Phonetics, Phonology, Prosody / Phonetics and Phoneme Processing / Phoneme Recognition Tasks     | Phoneme Segment Counting            | Librispeech-words (lib)   |
| Speech / Phonetics, Phonology, Prosody / Phonetics and Phoneme Processing / Phoneme Recognition Tasks     | Phone Segment Counting              | VoxAngeles (Chodroff et al., 2024)                              |
| Speech / Phonetics, Phonology, Prosody / Phonetics and Phoneme Processing / Phoneme Recognition Tasks     | SuperbPR                            | LibriSpeech-TestClean (Panayotov et al., 2015)                  |
| Speech / Phonetics, Phonology, Prosody / Phonetics and Phoneme Processing / Phoneme Recognition Tasks     | SuperbPR                            | LibriSpeech-TestOther (Panayotov et al., 2015)                  |
| Speech / Phonetics, Phonology, Prosody / Phonetics and Phoneme Processing / Phonological Feature Analysis | Phonological Feature Classification | VoxAngeles-ConsonantPlaceOfArticulation (Chodroff et al., 2024) |
| Speech / Phonetics, Phonology, Prosody / Phonetics and Phoneme Processing / Phonological Feature Analysis | Phonological Feature Classification | VoxAngeles-MannerOfArticulation (Chodroff et al., 2024)         |
| Speech / Phonetics, Phonology, Prosody / Phonetics and Phoneme Processing / Phonological Feature Analysis | Phonological Feature Classification | VoxAngeles-Phone (Chodroff et al., 2024)                        |
| Speech / Phonetics, Phonology, Prosody / Phonetics and Phoneme Processing / Phonological Feature Analysis | Phonological Feature Classification | VoxAngeles-VowelFrontness (Chodroff et al., 2024)               |

*Continued on next page*

| Domain   | Task  | Dataset  |
|--|---|--|
| Speech / Phonetics,<br>Phonology, Prosody /<br>Phonetics and Phoneme<br>Processing /<br>Phonological Feature<br>Analysis | Phonological Feature<br>Classification      | VoxAngeles-VowelHeight (Chodroff<br>et al., 2024)                                |
| Speech / Phonetics,<br>Phonology, Prosody /<br>Phonetics and Phoneme<br>Processing /<br>Phonological Feature<br>Analysis | Phonological Feature<br>Classification      | VoxAngeles-VowelRoundedness<br>(Chodroff et al., 2024)                           |
| Speech / Phonetics,<br>Phonology, Prosody /<br>Phonetics and Phoneme<br>Processing /<br>Pronunciation<br>Evaluation      | Heteronym<br>Differentiation                | HeteronymEn (self-constructed)   |
| Speech / Phonetics,<br>Phonology, Prosody /<br>Phonetics and Phoneme<br>Processing /<br>Pronunciation<br>Evaluation      | L2 English Accuracy                         | speechocean762-Ranking (Zhang<br>et al., 2021)                                   |
| Speech / Phonetics,<br>Phonology, Prosody /<br>Phonetics and Phoneme<br>Processing /<br>Pronunciation<br>Evaluation      | L2 English Accuracy                         | speechocean762-scoring (Zhang et al.,<br>2021)                                   |
| Speech / Phonetics,<br>Phonology, Prosody /<br>Phonetics and Phoneme<br>Processing /<br>Pronunciation<br>Evaluation      | Multilingual<br>Pronunciation<br>Similarity | VoxAngeles (Chodroff et al., 2024)   |
| Speech / Phonetics,<br>Phonology, Prosody /<br>Prosody / Prosodic<br>Feature Classification                              | Accent Classification                       | Accentdb Extended (Ahamad et al.,<br>2020)                                       |
| Speech / Phonetics,<br>Phonology, Prosody /<br>Prosody / Prosodic<br>Feature Classification                              | Stress Detection                            | MIRSD (CHEN & JANG, 2012;<br>Tseng, 2008; Tsai, 2015)                            |
| Speech / Phonetics,<br>Phonology, Prosody /<br>Prosody / Prosodic<br>Feature Classification                              | Third Tone Sandhi<br>Recognition            | NCCU Corpus of Spoken Taiwan<br>Mandarin (National Chengchi<br>University, n.d.) |

*Continued on next page*

| Domain   | Task                         | Dataset  |
|--|------------------------------|--|
| Speech / Phonetics, Phonology, Prosody / Prosody / Prosodic Quality Assessment | L2 English Fluency           | speechocean762-Ranking (Zhang et al., 2021)                          |
| Speech / Phonetics, Phonology, Prosody / Prosody / Prosodic Quality Assessment | L2 English Fluency           | speechocean762-Scoring (Zhang et al., 2021)                          |
| Speech / Phonetics, Phonology, Prosody / Prosody / Prosodic Quality Assessment | L2 English Prosodic          | speechocean762-Ranking (Zhang et al., 2021)                          |
| Speech / Phonetics, Phonology, Prosody / Prosody / Prosodic Quality Assessment | L2 English Prosodic          | speechocean762-Scoring (Zhang et al., 2021)                          |
| Speech / Phonetics, Phonology, Prosody / Prosody / Prosodic Quality Assessment | Prosody Naturalness          | ProsAudit-Lexical (de Seyssel et al., 2023)                          |
| Speech / Phonetics, Phonology, Prosody / Prosody / Prosodic Quality Assessment | Prosody Naturalness          | ProsAudit-Protosyntax (de Seyssel et al., 2023)                      |
| Speech / Safety & Security / Spoofing and Anti-Spoofing                        | Spoof Detection              | ASVspoof2015 (Wu et al., 2015)                                       |
| Speech / Safety & Security / Spoofing and Anti-Spoofing                        | Spoof Detection              | ASVspoof2017 (Delgado et al., 2018)                                  |
| Speech / Safety & Security / Synthetic Speech Detection                        | Deep Fake Voice Recognition  | DEEP-VOICE (Bird & Lotfi, 2023)                                      |
| Speech / Safety & Security / Synthetic Speech Detection                        | Enhancement Detection        | LibriTTS-TestClean (Zen et al., 2019)<br>WHAM (Wichern et al., 2019) |
| Speech / Safety & Security / Synthetic Speech Detection                        | Fraud Robocall Recognition   | CallHome (Katerenchuk et al., 2018)                                  |
| Speech / Safety & Security / Synthetic Speech Detection                        | Fraud Robocall Recognition   | Promo (self-constructed)   |
| Speech / Safety & Security / Synthetic Speech Detection                        | Fraud Robocall Recognition   | Robocall Prasad & Reaves   |
| Speech / Speaker & Language / Language / Language Identification               | HEAR Language Identification | VoxLingua107 Top10 (Valk & Alumäe, 2021)                             |

*Continued on next page*

| Domain   | Task                              | Dataset  |
|--|-----------------------------------|--|
| Speech / Speaker & Language / Language / Language Identification | Language Identification           | VoxForge (MacLean, 2018)   |
| Speech / Speaker & Language / Speaker / Speaker Characteristics  | Age Classification                | CommonVoiceCorpus-Test (Ardila et al., 2020)                           |
| Speech / Speaker & Language / Speaker / Speaker Characteristics  | Gender Recognition by Voice       | CommonVoice-DeltaSegment-15 (Ardila et al., 2020)                      |
| Speech / Speaker & Language / Speaker / Speaker Recognition      | HEAR Speaker Count Identification | LibriCount (Stöter et al., 2018)                                       |
| Speech / Speaker & Language / Speaker / Speaker Recognition      | Multi Speaker Detection           | VCTK (Yamagishi et al., 2019)  |
| Speech / Speaker & Language / Speaker / Speaker Recognition      | Speaker Counting                  | LibriTTS-TestClean (Zen et al., 2019)                                  |
| Speech / Speaker & Language / Speaker / Speaker Recognition      | Speaker Verification              | LibriSpeech-TestClean (Panayotov et al., 2015)                         |
| Speech / Speaker & Language / Speaker / Speaker Recognition      | Speaker Verification              | LibriSpeech-TestOther (Panayotov et al., 2015)                         |
| Speech / Speaker & Language / Speaker / Speaker Recognition      | Speaker Verification              | VCTK (Yamagishi et al., 2019)  |
| Speech / Speaker & Language / Speaker / Speaker Recognition      | SuperbSD                          | Libri2Mix-Test (Cosentino et al., 2020)                                |
| Speech / Speaker & Language / Speaker / Speaker Recognition      | SuperbSV                          | SuperbHiddenSet (self-reconstructed)                                   |
| Speech / Speech Enhancement / Degradation Detection              | Noise Detection                   | LJSpeech (Ito & Johnson, 2017)<br>MUSAN-Gaussian (Snyder et al., 2015) |
| Speech / Speech Enhancement / Degradation Detection              | Noise Detection                   | LJSpeech (Ito & Johnson, 2017)<br>MUSAN-Music (Snyder et al., 2015)    |
| Speech / Speech Enhancement / Degradation Detection              | Noise Detection                   | LJSpeech (Ito & Johnson, 2017)<br>MUSAN-Noise (Snyder et al., 2015)    |
| Speech / Speech Enhancement / Degradation Detection              | Noise Detection                   | LJSpeech (Ito & Johnson, 2017)<br>MUSAN-Speech (Snyder et al., 2015)   |
| Speech / Speech Enhancement / Degradation Detection              | Noise Detection                   | VCTK (Yamagishi et al., 2019)<br>MUSAN-Music (Snyder et al., 2015)     |

*Continued on next page*

| <b>Domain</b>                                       | <b>Task</b>                       | <b>Dataset</b>  |
|---|-----------------------------------|---|
| Speech / Speech Enhancement / Degradation Detection | Noise Detection                   | VCTK (Yamagishi et al., 2019)<br>MUSAN-Noise (Snyder et al., 2015)        |
| Speech / Speech Enhancement / Degradation Detection | Noise Detection                   | VCTK (Yamagishi et al., 2019)<br>MUSAN-Speech (Snyder et al., 2015)       |
| Speech / Speech Enhancement / Degradation Detection | Noise Detection                   | VCTK (Yamagishi et al., 2019)<br>MUSAN-Gaussian (Snyder et al., 2015)     |
| Speech / Speech Enhancement / Degradation Detection | Noise SNR Level Prediction        | VCTK (Yamagishi et al., 2019)<br>MUSAN-Gaussian (Snyder et al., 2015)     |
| Speech / Speech Enhancement / Degradation Detection | Noise SNR Level Prediction        | VCTK (Yamagishi et al., 2019)<br>MUSAN-Music (Snyder et al., 2015)        |
| Speech / Speech Enhancement / Degradation Detection | Noise SNR Level Prediction        | VCTK (Yamagishi et al., 2019)<br>MUSAN-Noise (Snyder et al., 2015)        |
| Speech / Speech Enhancement / Degradation Detection | Noise SNR Level Prediction        | VCTK (Yamagishi et al., 2019)<br>MUSAN-Speech (Snyder et al., 2015)       |
| Speech / Speech Enhancement / Degradation Detection | Reverberation Detection           | LJSpeech (Ito & Johnson, 2017)<br>RirsNoises-LargeRoom (Ko et al., 2017)  |
| Speech / Speech Enhancement / Degradation Detection | Reverberation Detection           | LJSpeech (Ito & Johnson, 2017)<br>RirsNoises-MediumRoom (Ko et al., 2017) |
| Speech / Speech Enhancement / Degradation Detection | Reverberation Detection           | LJSpeech (Ito & Johnson, 2017)<br>RirsNoises-SmallRoom (Ko et al., 2017)  |
| Speech / Speech Enhancement / Degradation Detection | Reverberation Detection           | VCTK (Yamagishi et al., 2019)<br>RirsNoises-LargeRoom (Ko et al., 2017)   |
| Speech / Speech Enhancement / Degradation Detection | Reverberation Detection           | VCTK (Yamagishi et al., 2019)<br>RirsNoises-MediumRoom (Ko et al., 2017)  |
| Speech / Speech Enhancement / Degradation Detection | Reverberation Detection           | VCTK (Yamagishi et al., 2019)<br>RirsNoises-SmallRoom (Ko et al., 2017)   |
| Speech / Speech Recognition / ASR Post-Processing   | N-Best Correction                 | Librispeech-TestOther (Panayotov et al., 2015)                            |
| Speech / Speech Recognition / Language              | AAVE Speech Recognition           | CORAAL (Kendall & Farrington, 2023)                                       |
| Speech / Speech Recognition / Language              | Code-switch Speech Recognition    | NTUML2021 (Yang et al., 2024a)  |
| Speech / Speech Recognition / Language              | Code-Switching Speech Recognition | ASCEND (Lovenia et al., 2022)   |

*Continued on next page*

| <b>Domain</b>  | <b>Task</b>                      | <b>Dataset</b>   |
|--|----------------------------------|--|
| Speech / Speech Recognition / Language                   | Multi-Lingual Speech Recognition | MLS-de (Pratap et al., 2020)                             |
| Speech / Speech Recognition / Language                   | Multi-Lingual Speech Recognition | MLS-en (Pratap et al., 2020)                             |
| Speech / Speech Recognition / Language                   | Multi-Lingual Speech Recognition | MLS-es (Pratap et al., 2020)                             |
| Speech / Speech Recognition / Language                   | Multi-Lingual Speech Recognition | MLS-fr (Pratap et al., 2020)                             |
| Speech / Speech Recognition / Language                   | Multi-Lingual Speech Recognition | MLS-it (Pratap et al., 2020)                             |
| Speech / Speech Recognition / Language                   | Multi-Lingual Speech Recognition | MLS-nl (Pratap et al., 2020)                             |
| Speech / Speech Recognition / Language                   | Multi-Lingual Speech Recognition | MLS-pl (Pratap et al., 2020)                             |
| Speech / Speech Recognition / Language                   | Multi-Lingual Speech Recognition | MLS-pt (Pratap et al., 2020)                             |
| Speech / Speech Recognition / Language                   | PTBR Speech Recognition          | CommonVoice17-Test (Ardila et al., 2020)                 |
| Speech / Speech Recognition / Language                   | SuperbASR                        | LibriSpeech-TestClean (Panayotov et al., 2015)           |
| Speech / Speech Recognition / Language                   | SuperbASR                        | LibriSpeech-TestOther (Panayotov et al., 2015)           |
| Speech / Speech Recognition / Language                   | SuperbOOD ASR Ar                 | CommonVoice7-Test (Ardila et al., 2020)                  |
| Speech / Speech Recognition / Language                   | SuperbOOD ASR Es                 | CommonVoice7-Test (Ardila et al., 2020)                  |
| Speech / Speech Recognition / Language                   | SuperbOOD ASR Spon               | CHIME6-Test (Barker et al., 2018; Watanabe et al., 2020) |
| Speech / Speech Recognition / Language                   | SuperbOOD ASr Zh                 | CommonVoice7-Test (Ardila et al., 2020)                  |
| Speech / Speech Recognition / Speaker                    | Target Speaker ASR               | AMI-test (Carletta et al., 2005)                         |
| Speech / Speech Recognition / Specific Recognition Tasks | Multi-Speaker Detection          | LibriSpeech-TestClean (Panayotov et al., 2015)           |
| Speech / Speech Recognition / Specific Recognition Tasks | Speech Command Recognition       | AudioMNIST (Becker et al., 2024)                         |
| Speech / Speech Recognition / Specific Recognition Tasks | Speech Text Matching             | LibriSpeech-TestClean (Panayotov et al., 2015)           |
| Speech / Speech Recognition / Specific Recognition Tasks | Speech Text Matching             | LibriSpeech-TestOther (Panayotov et al., 2015)           |

*Continued on next page*



| Domain   | Task                          | Dataset  |
|--|-------------------------------|--|
| Speech / Speech Recognition / Specific Recognition Tasks                         | Speech Text Matching          | LJSpeech (Ito & Johnson, 2017)                 |
| Speech / Speech Recognition / Specific Recognition Tasks                         | Spoken Term Detection         | LibriSpeech-TestClean (Panayotov et al., 2015) |
| Speech / Speech Recognition / Specific Recognition Tasks                         | Spoken Term Detection         | LibriSpeech-TestOther (Panayotov et al., 2015) |
| Speech / Speech Recognition / Specific Recognition Tasks                         | Spoken Term Detection         | LJSpeech (Ito & Johnson, 2017)                 |
| Speech / Speech Recognition / Specific Recognition Tasks                         | SuperbKS                      | SpeechCommandsV1-Test (Warden, 2018)           |
| Speech / Speech Recognition / Specific Recognition Tasks                         | SuperbQbE                     | Quesst14-Eval (Anguera et al., 2015)           |
| Speech / Speech, Voice, Hearing Disorder / Disorder Detection and Classification | Stuttering Detection          | SEP-28k (Lea et al., 2021)                     |
| Speech / Speech, Voice, Hearing Disorder / Disorder Detection and Classification | Voice Disorder Classification | VOICED (Goldberger et al., 2000)               |
| Speech / Spoken Language Understanding / Intent & Meaning                        | Conversation Matching         | EnShortConversation (CEEC)                     |
| Speech / Spoken Language Understanding / Intent & Meaning                        | Dialogue Act Classification   | SLUE-HVB (Shon et al., 2022)                   |
| Speech / Spoken Language Understanding / Intent & Meaning                        | Dialogue Act Classification   | DailyTalk (Lee et al., 2023)                   |
| Speech / Spoken Language Understanding / Intent & Meaning                        | Dialogue Act Pairing          | DailyTalk (Lee et al., 2023)                   |
| Speech / Spoken Language Understanding / Intent & Meaning                        | SuperbIC                      | SLURP (Bastianelli et al., 2020)               |

*Continued on next page*

| <b>Domain</b>   | <b>Task</b>   | <b>Dataset</b>                          |
|---|---|---|
| Speech / Spoken<br>Language<br>Understanding / Intent<br>& Meaning                  | Intent Classification<br>(Action)                                 | SLURP-Action (Bastianelli et al., 2020) |
| Speech / Spoken<br>Language<br>Understanding / Intent<br>& Meaning                  | Intent Classification<br>(Intent)                                 | SLURP-Intent (Bastianelli et al., 2020) |
| Speech / Spoken<br>Language<br>Understanding / Intent<br>& Meaning                  | Named Entity<br>Localization                                      | SLUE-VoxPopuli (Shon et al., 2022)      |
| Speech / Spoken<br>Language<br>Understanding / Intent<br>& Meaning                  | Named Entity<br>Recognition                                       | SLUE-VoxPopuli (Shon et al., 2022)      |
| Speech / Spoken<br>Language<br>Understanding / Intent<br>& Meaning                  | Sarcasm Detection   | Mustard (Castro et al., 2019)           |
| Speech / Spoken<br>Language<br>Understanding / Intent<br>& Meaning                  | Semantic Textual<br>Similarity                                    | SpokenSTS (Agirre et al., 2012)         |
| Speech / Spoken<br>Language<br>Understanding / Intent<br>& Meaning                  | Speech Sentiment<br>Analysis                                      | MELD (Poria et al., 2019)               |
| Speech / Spoken<br>Language<br>Understanding / Intent<br>& Meaning                  | Spoken Digit<br>Arithmetic  | AudioMNIST (Becker et al., 2024)        |
| Speech / Spoken<br>Language<br>Understanding / Intent<br>& Meaning                  | SuperbSF  | AudioSnips-Test (Yang et al., 2021)     |
| Speech / Spoken<br>Language<br>Understanding /<br>Linguistic Structure &<br>Grammar | Code-Switching<br>Semantic Grammar<br>Acceptability<br>Comparison | CSZS-zh-en (Huang et al., 2024b)        |
| Speech / Spoken<br>Language<br>Understanding /<br>Linguistic Structure &<br>Grammar | Nonce Word<br>Detection   | sWUGGY (Dunbar et al.)                  |

*Continued on next page*

| Domain  | Task  | Dataset                          |
|---|---|----------------------------------|
| Speech / Spoken Language Understanding / Linguistic Structure & Grammar | PoS Estimation<br>LibriTTS PoS                    | LibriTTS (Zen et al., 2019)      |
| Speech / Spoken Language Understanding / Linguistic Structure & Grammar | PoS Estimation<br>LibriTTS PoS with transcription | LibriTTS (Zen et al., 2019)      |
| Speech / Spoken Language Understanding / Linguistic Structure & Grammar | Sentence Grammar Acceptability                    | sBLIMP (Dunbar et al.)           |
| Speech / Spoken Language Understanding / Speech Translation             | SuperbST  | CoVoST2-Test (Wang et al., 2021) |

## C TASK-LEVEL EVALUATION RESULTS

Table 5: Task-level Evaluation Results: GPT4o

| Task   | Metric           | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM  | MU-LLaMA | GAMA-IT |
|--|------------------|---------------|---------|------------|-------------|-----------------|----------------------|---------|----------|---------|
| Chord Classification   | LLM-C $\uparrow$ | 4.00%         | 10.00%  | 0.50%      | 9.00%       | 44.00%          | 28.50%               | 44.50%  | 51.00%   | 0.00%   |
| MARBLE Key Detection (GiantSteps - Key)                            | LLM-C $\uparrow$ | 3.00%         | 0.50%   | 1.50%      | 1.00%       | 3.00%           | 17.00%               | 2.00%   | 0.50%    | 0.00%   |
| HEAR Music Transcription (MAESTRO)                                 | NAR $\downarrow$ | 92.43%        | 100.00% | 99.46%     | 100.00%     | 96.76%          | 96.76%               | 96.76%  | 100.00%  | 100.00% |
| HEAR Music Transcription (MAESTRO)                                 | TER $\downarrow$ | 1.0065        | -       | 36.0000    | -           | 8.3438          | 1.6875               | 6.25    | -        | -       |
| HEAR Percussion Instruments Tonic Classification (Mridangam Tonic) | LLM-C $\uparrow$ | 2.00%         | 2.50%   | 12.70%     | 13.40%      | 11.30%          | 11.80%               | 12.80%  | 12.50%   | 5.70%   |
| Instrument Pitch Classification (Nsynth)                           | LLM-C $\uparrow$ | 0.56%         | 0.56%   | 0.11%      | 0.22%       | 0.67%           | 18.11%               | 0.78%   | 0.33%    | 0.78%   |
| Pitch Extraction By Lyrics (CSD)                                   | NAR $\downarrow$ | 37.50%        | 94.00%  | 47.50%     | 73.50%      | 65.50%          | 84.50%               | 76.00%  | 98.00%   | 68.50%  |
| Pitch Extraction By Lyrics (CSD)                                   | TER $\downarrow$ | 146.28%       | 116.67% | 226.21%    | 639.24%     | 106.45%         | 165.96%              | 136.11% | 100.00%  | 195.18% |

*Continued on next page*

| Task  | Metric           | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM  | MU-LLaMA | GAMA-IT |
|---|------------------|---------------|---------|------------|-------------|-----------------|----------------------|---------|----------|---------|
| Emotion Classification In Songs (EMOTIFY)                             | LLM-C $\uparrow$ | 10.00%        | 0.00%   | 1.67%      | 3.33%       | 13.33%          | 8.33%                | 8.33%   | 1.67%    | 0.00%   |
| MARBLE Emotion Detection (MTG)  | LLM-C $\uparrow$ | 0.00%         | 1.50%   | 0.00%      | 0.00%       | 2.70%           | 0.80%                | 0.80%   | 0.80%    | 0.10%   |
| HEAR Music Genre Classification (ISMIR04)                             | LLM-C $\uparrow$ | 15.58%        | 7.04%   | 31.16%     | 44.22%      | 33.67%          | 37.69%               | 16.58%  | 27.64%   | 2.51%   |
| MARBLE Genre Classification (MTG)                                     | LLM-C $\uparrow$ | 0.00%         | 1.60%   | 0.10%      | 0.20%       | 3.20%           | 0.10%                | 0.20%   | 0.20%    | 0.00%   |
| MARBLE Music Tagging (MagnaTagATune)                                  | LLM-C $\uparrow$ | 5.50%         | 1.00%   | 3.50%      | 9.00%       | 4.50%           | 1.50%                | 1.00%   | 3.50%    | 7.50%   |
| MARBLE Music Tagging (MTG)  | LLM-C $\uparrow$ | 0.10%         | 0.30%   | 0.00%      | 0.00%       | 2.90%           | 0.30%                | 0.10%   | 0.50%    | 0.10%   |
| Music Genre Classification (FMA)                                      | LLM-C $\uparrow$ | 10.71%        | 1.79%   | 9.82%      | 16.07%      | 8.04%           | 15.18%               | 9.82%   | 9.82%    | 4.46%   |
| HEAR Percussion Instruments Classification (Beijing Opera Percussion) | LLM-C $\uparrow$ | 3.39%         | 4.24%   | 10.17%     | 20.76%      | 16.10%          | 27.12%               | 24.15%  | 22.03%   | 11.44%  |
| HEAR Percussion Instruments Stroke Classification (Mridangam Stroke)  | LLM-C $\uparrow$ | 0.10%         | 3.80%   | 7.60%      | 4.90%       | 11.50%          | 10.80%               | 13.70%  | 10.30%   | 0.10%   |
| Instrument Classification (Nsynth)                                    | LLM-C $\uparrow$ | 1.23%         | 9.56%   | 13.97%     | 14.58%      | 17.89%          | 54.41%               | 11.89%  | 6.99%    | 10.17%  |
| Instrument Combination Recognition (OpenMIC-2018)                     | LLM-C $\uparrow$ | 4.20%         | 13.51%  | 17.72%     | 24.02%      | 25.83%          | 24.92%               | 20.42%  | 19.52%   | 8.11%   |
| Instrument Source Classification (Nsynth)                             | LLM-C $\uparrow$ | 7.22%         | 28.33%  | 37.33%     | 43.56%      | 32.67%          | 35.22%               | 30.89%  | 38.00%   | 35.89%  |
| MARBLE Instrument Classification (MTG)                                | LLM-C $\uparrow$ | 0.40%         | 2.30%   | 0.00%      | 0.00%       | 6.40%           | 1.40%                | 0.20%   | 0.60%    | 0.20%   |
| Singing Automatic MOS Prediction (SingMOS)                            | NAR $\downarrow$ | 99.50%        | 67.72%  | 0.33%      | 11.15%      | 17.30%          | 0.00%                | 22.96%  | 97.34%   | 98.34%  |
| Singing Automatic MOS Prediction (SingMOS)                            | MSE $\downarrow$ | 12.17         | 2.3872  | 10.0777    | 11.2308     | 14.8895         | 1.4485               | 11.421  | 4.53     | 2.08    |
| Singing Automatic MOS Prediction (SingMOS)                            | KTAU $\uparrow$  | -0.3333       | 0.0126  | -0.1105    | 0.0166      | 0.0712          | 0.0099               | 0.0103  | 0.1705   | 0.3539  |
| Singing Automatic MOS Prediction (SingMOS)                            | LCC $\uparrow$   | -0.6747       | -0.0706 | -0.1354    | 0.036       | 0.0824          | 0.0255               | -0.0216 | 0.1866   | 0.3846  |
| Singing Automatic MOS Prediction (SingMOS)                            | SRCC $\uparrow$  | -0.5          | 0.0183  | -0.1325    | 0.0231      | 0.0872          | 0.0124               | 0.0118  | 0.2281   | 0.4247  |

Continued on next page

| Task  | Metric      | Whisper-LLaMA | LTU-AS   | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM    | MU-LLaMA | GAMA-IT  |
|---|-------------|---------------|----------|------------|-------------|-----------------|----------------------|-----------|----------|----------|
| MARBLE Beat Tracking (ASAP)                                     | NAR↓        | 100.00%       | 100.00%  | 100.00%    | 100.00%     | 100.00%         | 100.00%              | 100.00%   | 100.00%  | 100.00%  |
| MARBLE Beat Tracking (ASAP)                                     | ERR↓        | -             | -        | -          | -           | -               | -                    | -         | -        | -        |
| Music Beat Tracking (ASAP)                                      | NAR↓        | 54.55%        | 74.55%   | 51.82%     | 63.64%      | 62.73%          | 40.91%               | 47.27%    | 80.00%   | 84.55%   |
| Music Beat Tracking (ASAP)                                      | Miss Time↓  | 135.6077      | 183.6786 | 118.4047   | 77.9687     | 157.5917        | 76.0904              | 111.7759  | 240      | 93.8751  |
| Audio Editing Identification (People's Speech)                  | NAR↓        | 84.38%        | 87.50%   | 62.50%     | 75.00%      | 75.00%          | 68.75%               | 96.88%    | 65.62%   | 37.50%   |
| Audio Editing Identification (People's Speech)                  | ACC↑        | 204.3094      | 53.6938  | 165.0617   | 43.1693     | 184.695         | 154.6202             | 100       | 164.0558 | 154.8122 |
| Scene Fake Detection (ASPIRE)                                   | LLM-C↑      | 22.00%        | 46.00%   | 48.00%     | 37.00%      | 49.00%          | 54.00%               | 49.00%    | 20.00%   | 41.00%   |
| Audio Deep Fake Detection (LJSpeech, WaveFake, MUSDB18HQ)       | LLM-C↑      | 8.25%         | 24.56%   | 22.79%     | 38.31%      | 31.43%          | 27.90%               | 30.84%    | 39.10%   | 25.54%   |
| Singing Voice Deepfake Detection (CtrSVDD, ACEKiSing, M4Singer) | LLM-C↑      | 5.56%         | 22.06%   | 22.06%     | 26.19%      | 24.44%          | 21.43%               | 9.21%     | 20.95%   | 9.84%    |
| Audio Duration Prediction (NTUML2021)                           | NAR↓        | 33.00%        | 37.00%   | 10.50%     | 55.00%      | 0.00%           | 0.00%                | 0.50%     | 46.50%   | 14.00%   |
| Audio Duration Prediction (NTUML2021)                           | MSE↓        | 13.3657       | 39.1111  | 31652.7877 | 2216.7111   | 28.985          | 62.515               | 3581.0352 | 1527.785 | 62.7674  |
| HEAR Music Speech Classification (MAESTRO, Librispeech)         | LLM-C↑      | 72.50%        | 77.50%   | 91.50%     | 99.50%      | 98.00%          | 89.50%               | 59.50%    | 51.50%   | 88.00%   |
| Sound Effect Detection (RemFx)                                  | LLM-C↑      | 2.00%         | 6.50%    | 14.17%     | 17.67%      | 15.83%          | 19.83%               | 17.50%    | 17.00%   | 15.83%   |
| Speech Detection (LibriSpeech-TestClean)                        | LLM-C↑      | 48.00%        | 29.00%   | 46.00%     | 66.00%      | 45.50%          | 56.50%               | 35.00%    | 46.50%   | 38.50%   |
| Speech Detection (LibriSpeech-TestOther)                        | LLM-C↑      | 52.00%        | 31.50%   | 46.00%     | 66.50%      | 44.50%          | 55.50%               | 37.00%    | 45.50%   | 36.00%   |
| Speech Detection (LJSpeech)                                     | LLM-C↑      | 54.00%        | 44.50%   | 60.50%     | 74.00%      | 57.00%          | 45.00%               | 53.00%    | 49.00%   | 55.50%   |
| Children Song Transcript Verification (CSD)                     | WER↓        | 58.68%        | 102.68%  | 104.02%    | 100.99%     | 128.67%         | 94.93%               | 155.48%   | 100.00%  | 100.12%  |
| Lyric Translation (SingSet)                                     | Sacre Bleu↑ | 2.1093        | 1.0357   | 2.8126     | 0.9941      | 3.3975          | 3.9155               | 1.8969    | 0.1593   | 0.0739   |
| Song Lyric Recognition (SingSet)                                | MER↓        | 101.93%       | 164.53%  | 356.36%    | 266.00%     | 118.41%         | 115.62%              | 481.13%   | 117.53%  | 168.52%  |

Continued on next page

| Task   | Metric          | Whisper-LLaMA | LFU-AS | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM | MU-LLaMA | GAMA-IT |
|--|-----------------|---------------|--------|------------|-------------|-----------------|----------------------|--------|----------|---------|
| MARBLE Vocal Technique Detection (VocalSet)                          | LLM-C↑          | 8.00%         | 1.00%  | 11.00%     | 6.50%       | 6.50%           | 15.50%               | 8.00%  | 8.50%    | 2.50%   |
| Audio Segment Retrieval (Clotho)                                     | NAR↓            | 80.00%        | 78.00% | 2.00%      | 4.00%       | 12.00%          | 4.00%                | 14.00% | 82.00%   | 6.00%   |
| Audio Segment Retrieval (Clotho)                                     | IoU↑            | 8.17%         | 22.17% | 8.25%      | 8.68%       | 4.58%           | 7.07%                | 6.90%  | 0.00%    | 10.45%  |
| HEAR Sound Event Detection (DCASE2016 Task2)                         | LLM-C↑          | 0.00%         | 0.00%  | 21.43%     | 35.71%      | 7.14%           | 14.29%               | 21.43% | 21.43%   | 7.14%   |
| Multi-channel Sound Event Understanding (STARSS23)                   | LLM-C↑          | 38.73%        | 4.23%  | 37.32%     | 26.06%      | 38.03%          | 41.55%               | 21.83% | 4.93%    | 11.27%  |
| Animal Classification (WaveSource-Test)                              | LLM-C↑          | 9.25%         | 40.50% | 58.25%     | 67.00%      | 75.75%          | 34.00%               | 12.00% | 4.25%    | 15.25%  |
| Bird Sound Detection (Warblr10k)                                     | LLM-C↑          | 28.50%        | 43.00% | 74.50%     | 75.00%      | 78.50%          | 79.00%               | 75.00% | 40.50%   | 33.00%  |
| Cat Emotion Classification (Cat Sound Classification Dataset V2)     | LLM-C↑          | 6.00%         | 0.00%  | 2.00%      | 4.00%       | 6.00%           | 14.00%               | 8.00%  | 2.00%    | 14.00%  |
| Cornell Birdcall Identification                                      | LLM-C↑          | 0.00%         | 0.00%  | 0.00%      | 3.33%       | 10.00%          | 10.00%               | 13.33% | 3.33%    | 0.00%   |
| Environmental Sound Classification (ESC50-Animals)                   | LLM-C↑          | 14.50%        | 8.00%  | 14.00%     | 3.00%       | 82.00%          | 36.00%               | 5.00%  | 0.00%    | 2.00%   |
| HEAR Beehive States Classification (Beehive States)                  | LLM-C↑          | 43.00%        | 35.00% | 54.00%     | 42.50%      | 19.00%          | 42.50%               | 37.00% | 37.50%   | 19.00%  |
| Domestic Environment Sound Event Detection (DESED-PublicEval)        | NAR↓            | 91.39%        | 99.44% | 99.72%     | 99.72%      | 99.17%          | 99.72%               | 99.72% | 100.00%  | 99.72%  |
| Domestic Environment Sound Event Detection (DESED-PublicEval)        | Event-based F1↑ | 0             | 0.6    | 0          | 0           | 0               | 0.01                 | 0      | 0        | 0       |
| Emergency Traffic Detection (Large-Scale-Audio-dataset)              | LLM-C↑          | 39.00%        | 34.00% | 55.50%     | 7.25%       | 65.50%          | 74.50%               | 17.50% | 43.50%   | 15.83%  |
| Environmental Sound Classification (ESC50-Exterior And Urban Noises) | LLM-C↑          | 4.50%         | 14.50% | 0.50%      | 0.50%       | 77.00%          | 25.00%               | 6.50%  | 0.00%    | 6.50%   |

*Continued on next page*

| Task  | Metric      | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM | MU-LLaMA | GAMA-IT |
|---|-------------|---------------|---------|------------|-------------|-----------------|----------------------|--------|----------|---------|
| Environmental Sound Classification (ESC50-Interior And Domestic Sounds)         | LLM-C↑      | 4.50%         | 3.50%   | 0.50%      | 0.00%       | 59.00%          | 2.50%                | 4.00%  | 0.00%    | 0.00%   |
| Environmental Sound Classification (ESC50-Natural Soundscapes And Water Sounds) | LLM-C↑      | 3.00%         | 7.50%   | 6.50%      | 1.00%       | 78.00%          | 5.00%                | 7.00%  | 0.00%    | 7.00%   |
| Environmental Sound Classification (UrbanSound8K-Urban Noises)                  | LLM-C↑      | 10.00%        | 3.26%   | 1.40%      | 0.00%       | 39.53%          | 4.88%                | 10.47% | 0.00%    | 0.70%   |
| Environment Recognition (ESC50)   | LLM-C↑      | 7.02%         | 38.60%  | 42.11%     | 19.30%      | 56.14%          | 63.16%               | 40.35% | 10.53%   | 31.58%  |
| HEAR Environmental Sound Classification (ESC50)                                 | LLM-C↑      | 4.80%         | 25.30%  | 61.00%     | 73.40%      | 68.80%          | 43.90%               | 9.80%  | 0.40%    | 35.30%  |
| HEAR Vocal Imitation Classification (Vocal Imitations)                          | LLM-C↑      | 15.17%        | 5.50%   | 1.17%      | 1.83%       | 18.67%          | 11.17%               | 17.50% | 3.83%    | 5.83%   |
| Audio Spatial Distance Prediction (Spatial LibriSpeech)                         | NAR↓        | 50.50%        | 59.00%  | 62.50%     | 100.00%     | 85.50%          | 58.00%               | 44.00% | 97.00%   | 71.00%  |
| Audio Spatial Distance Prediction (Spatial LibriSpeech)                         | MEDAE↓      | 0.6651        | 0.933   | 0.449      | -           | 1.5011          | 0.6915               | 0.9011 | 0.9212   | 0.5     |
| How Far Are You (3DSpeaker)   | LLM-C↑      | 26.00%        | 13.00%  | 4.50%      | 8.00%       | 30.00%          | 26.00%               | 32.00% | 18.00%   | 2.50%   |
| Sound Position Prediction   | NAR↓        | 81.25%        | 100.00% | 100.00%    | 100.00%     | 100.00%         | 25.00%               | 93.75% | 100.00%  | 100.00% |
| Sound Position Prediction   | Angle Diff↓ | 1.4006        | -       | -          | -           | -               | 1.1063               | 0.8681 | -        | -       |
| Dialogue Emotion Classification (Daily Talk)                                    | LLM-C↑      | 33.50%        | 16.50%  | 27.00%     | 11.50%      | 33.00%          | 54.00%               | 18.50% | 19.00%   | 33.00%  |
| Emoji Grounded Speech Emotion Recognition (RAVDESS)                             | LLM-C↑      | 1.10%         | 0.00%   | 0.00%      | 0.10%       | 1.30%           | 1.20%                | 0.10%  | 1.90%    | 48.50%  |
| Emotion Change Detection (Ravdess)  | LLM-C↑      | 0.08%         | 0.42%   | 0.58%      | 0.96%       | 0.04%           | 2.96%                | 0.08%  | 0.58%    | 0.04%   |
| Emotion Recognition (MELD)  | LLM-C↑      | 22.50%        | 14.50%  | 19.00%     | 9.50%       | 43.50%          | 36.50%               | 40.50% | 0.50%    | 2.50%   |
| HEAR Emotion Recognition (CREMAD)   | LLM-C↑      | 6.70%         | 9.30%   | 18.80%     | 18.10%      | 62.50%          | 61.10%               | 25.60% | 12.40%   | 13.20%  |
| SuperbER (RAVDESS)  | LLM-C↑      | 12.50%        | 25.83%  | 12.08%     | 12.50%      | 70.42%          | 75.83%               | 12.08% | 12.92%   | 0.00%   |

Continued on next page

| Task   | Metric    | Whisper-LLaMA | LTU-AS    | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM     | MU-LLaMA | GAMA-IT |
|--|-----------|---------------|-----------|------------|-------------|-----------------|----------------------|------------|----------|---------|
| Covid19 Cough Audio Classification (CoughVid)                                    | LLM-C↑    | 0.93%         | 1.03%     | 0.00%      | 0.00%       | 4.12%           | 0.72%                | 4.22%      | 21.91%   | 0.21%   |
| Environmental Sound Classification (ESC50-Human And Non Speech Sounds)           | LLM-C↑    | 17.50%        | 11.00%    | 8.50%      | 3.00%       | 82.00%          | 29.00%               | 16.50%     | 0.00%    | 2.00%   |
| Human Non Speech Sound Recognition (Nonspeech7k, CommonVoice)                    | LLM-C↑    | 2.86%         | 5.00%     | 44.29%     | 30.71%      | 27.14%          | 23.57%               | 14.29%     | 17.14%   | 32.86%  |
| Human Screaming Detection (Environmentdb)  | LLM-C↑    | 52.50%        | 25.00%    | 52.50%     | 65.00%      | 62.50%          | 87.50%               | 47.50%     | 42.50%   | 60.00%  |
| Vocal Sound Recognition (VocalSound)   | LLM-C↑    | 29.86%        | 1.53%     | 11.53%     | 4.31%       | 75.56%          | 30.97%               | 14.44%     | 0.28%    | 2.78%   |
| Phoneme Segment Counting (Librispeech-words)                                     | NAR↓      | 16.29%        | 48.64%    | 39.09%     | 11.62%      | 5.99%           | 1.73%                | 0.50%      | 2.67%    | 72.22%  |
| Phoneme Segment Counting (Librispeech-words)                                     | ACC↑      | 13.56%        | 11.77%    | 11.84%     | 8.25%       | 5.42%           | 13.46%               | 9.55%      | 0.47%    | 2.29%   |
| Phoneme Segment Counting (Librispeech-words)                                     | Abs Diff↓ | 2.5378        | 1012.8103 | 6.5285     | 50.7724     | 7.1857          | 3.0622               | 7.9109     | 15.6777  | 5.0491  |
| Phone Segment Counting (VoxAngeles)  | NAR↓      | 19.19%        | 57.24%    | 41.88%     | 12.56%      | 0.41%           | 16.11%               | 15.04%     | 41.47%   | 18.79%  |
| Phone Segment Counting (VoxAngeles)  | ACC↑      | 13.20%        | 1.69%     | 15.81%     | 20.04%      | 1.77%           | 6.60%                | 10.11%     | 9.14%    | 3.80%   |
| Phone Segment Counting (VoxAngeles)  | Abs Diff↓ | 41.236        | 3759.3878 | 5.4857     | 4.9331      | 3.1286          | 2.8706               | 45754.9952 | 5.4242   | 2.8694  |
| SuperbPR (LibriSpeech-TestClean)   | PER↓      | 100.12%       | 102.75%   | 25.36%     | 24.60%      | 100.58%         | 100.96%              | 99.99%     | 110.27%  | 100.37% |
| SuperbPR (LibriSpeech-TestOther)   | PER↓      | 100.17%       | 101.08%   | 22.79%     | 22.91%      | 100.62%         | 100.88%              | 99.99%     | 111.18%  | 100.61% |
| Phonological Feature Classification (VoxAngeles-Consonant Place Of Articulation) | LLM-C↑    | 25.77%        | 1.37%     | 1.88%      | 1.37%       | 1.54%           | 3.24%                | 1.88%      | 1.54%    | 0.51%   |
| Phonological Feature Classification (VoxAngeles-Manner Of Articulation)          | LLM-C↑    | 17.42%        | 9.36%     | 2.06%      | 1.50%       | 5.62%           | 6.74%                | 6.84%      | 7.96%    | 8.99%   |
| Phonological Feature Classification (VoxAngeles-Phone)                           | LLM-C↑    | 5.73%         | 0.09%     | 1.76%      | -           | 3.60%           | 3.97%                | 0.00%      | 0.00%    | 0.09%   |

Continued on next page



| Task  | Metric | Whisper-LLaMA | LTU-AS | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM | MU-LLaMA | GAMA-IT |
|---|--------|---------------|--------|------------|-------------|-----------------|----------------------|--------|----------|---------|
| Phonological Feature Classification (VoxAngeles-Vowel Frontness)      | LLM-C↑ | 50.71%        | 17.52% | 41.14%     | 38.09%      | 48.47%          | 57.23%               | 36.46% | 42.97%   | 10.39%  |
| Phonological Feature Classification (VoxAngeles-Vowel Height)         | LLM-C↑ | 24.44%        | 17.11% | 31.77%     | 35.23%      | 29.53%          | 36.86%               | 37.88% | 39.10%   | 24.44%  |
| Phonological Feature Classification (VoxAngeles-Vowel Roundedness)    | LLM-C↑ | 43.38%        | 21.18% | 27.29%     | 61.30%      | 46.44%          | 69.65%               | 38.90% | 21.79%   | 18.94%  |
| Heteronym Differentiation (HeteronymEn)                               | LLM-C↑ | 55.00%        | 26.00% | 37.00%     | 45.00%      | 52.00%          | 44.00%               | 51.00% | 30.00%   | 20.00%  |
| L2 English Accuracy (speechocean762-Ranking)                          | LLM-C↑ | 24.72%        | 33.89% | 50.00%     | 49.72%      | 34.44%          | 50.00%               | 41.67% | 48.06%   | 45.56%  |
| L2 English Accuracy (speechocean762-scoring)                          | NAR↓   | 15.44%        | 95.03% | 0.40%      | 0.27%       | 37.72%          | 0.94%                | 9.53%  | 90.47%   | 50.07%  |
| L2 English Accuracy (speechocean762-scoring)                          | PCC↑   | 0.0185        | -0.183 | 0.0633     | 0.0438      | 0.0293          | -0.0159              | 0.0439 | 0.0727   | 0.0151  |
| Multi-lingual Pronunciation Similarity (VoxAngeles)                   | LLM-C↑ | 38.40%        | 13.50% | 47.20%     | 48.50%      | 23.70%          | 44.70%               | 48.00% | 25.10%   | 40.60%  |
| Accent Classification (Accentdb Extended)                             | LLM-C↑ | 17.50%        | 5.50%  | 3.00%      | 4.50%       | 26.50%          | 14.00%               | 7.00%  | 27.00%   | 4.00%   |
| Stress Detection (MIRSD)  | LLM-C↑ | 15.50%        | 3.00%  | 2.00%      | 13.00%      | 16.50%          | 23.50%               | 25.50% | 0.00%    | 1.00%   |
| Third Tone Sandhi Recognition (NCCU Corpus of Spoken Taiwan Mandarin) | NAR↓   | 59.38%        | 93.75% | 100.00%    | 87.50%      | 84.38%          | 84.38%               | 28.12% | 96.88%   | 100.00% |
| Third Tone Sandhi Recognition (NCCU Corpus of Spoken Taiwan Mandarin) | IoU↑   | 0.2179        | 0.5    | 0          | 0.75        | 0               | 0                    | 0.1304 | 0        | 0       |
| L2 English Fluency (speechocean762-Ranking)                           | LLM-C↑ | 31.39%        | 21.94% | 50.00%     | 50.00%      | 40.28%          | 50.56%               | 27.78% | 49.44%   | 45.00%  |
| L2 English Fluency (speechocean762-Scoring)                           | NAR↓   | 14.67%        | 88.16% | 0.00%      | 0.67%       | 8.75%           | 0.81%                | 41.86% | 63.53%   | 55.72%  |
| L2 English Fluency (speechocean762-Scoring)                           | PCC↑   | 0.0055        | 0.0332 | 0.0292     | 0.0183      | -0.0422         | -0.0858              | 0.0222 | 0.0505   | 0.0532  |

*Continued on next page*

| Task   | Metric | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM | MU-LLaMA | GAMA-IT |
|--|--------|---------------|---------|------------|-------------|-----------------|----------------------|--------|----------|---------|
| L2 English Prosodic (speechoccean762-Ranking)    | LLM-C↑ | 26.39%        | 35.00%  | 46.67%     | 52.50%      | 34.44%          | 51.11%               | 34.44% | 49.72%   | 42.22%  |
| L2 English Prosodic (speechoccean762-Scoring)    | NAR↓   | 39.43%        | 94.75%  | 0.27%      | 1.21%       | 71.47%          | 3.10%                | 11.98% | 94.62%   | 82.91%  |
| L2 English Prosodic (speechoccean762-Scoring)    | PCC↑   | 0.0435        | -0.1747 | 0.0427     | 0.0775      | 0.1446          | 0.0201               | 0.0461 | -0.1535  | 0.0973  |
| Prosody Naturalness (ProsAudit-Lexical)          | LLM-C↑ | 49.81%        | 32.82%  | 48.26%     | 47.10%      | 21.62%          | 51.74%               | 54.83% | 47.10%   | 5.41%   |
| Prosody Naturalness (ProsAudit-Protosyntax)      | LLM-C↑ | 53.44%        | 31.68%  | 46.56%     | 46.18%      | 25.95%          | 50.76%               | 43.51% | 49.24%   | 4.20%   |
| Spoof Detection (ASVspoof2015)                   | LLM-C↑ | 41.00%        | 19.00%  | 55.00%     | 14.00%      | 13.00%          | 23.00%               | 19.50% | 69.50%   | 0.50%   |
| Spoof Detection (ASVspoof2017)                   | LLM-C↑ | 49.00%        | 4.50%   | 63.00%     | 38.00%      | 27.00%          | 32.50%               | 21.00% | 64.00%   | 0.00%   |
| Deep Fake Voice Recognition (DEEP-VOICE)         | LLM-C↑ | 27.25%        | 7.75%   | 50.75%     | 40.00%      | 48.75%          | 50.50%               | 43.50% | 45.00%   | 19.75%  |
| Enhancement Detection (LibriTTS-TestClean, WHAM) | LLM-C↑ | 50.50%        | 30.50%  | 30.00%     | 46.50%      | 55.00%          | 56.50%               | 30.00% | 50.00%   | 21.00%  |
| Fraud Robocall Recognition (CallHome)            | LLM-C↑ | 73.33%        | 13.33%  | 0.00%      | 100.00%     | 100.00%         | 100.00%              | 90.00% | 26.67%   | 100.00% |
| Fraud Robocall Recognition (Promo)               | LLM-C↑ | 47.37%        | 10.53%  | 0.00%      | 73.68%      | 57.89%          | 63.16%               | 63.16% | 5.26%    | 100.00% |
| Fraud Robocall Recognition (Robocall)            | LLM-C↑ | 94.87%        | 51.28%  | 100.00%    | 20.51%      | 51.28%          | 61.54%               | 30.77% | 30.77%   | 0.00%   |
| HEAR Language Identification (VoxLingua107)      | LLM-C↑ | 92.18%        | 1.95%   | 6.28%      | 11.01%      | 18.00%          | 36.11%               | 1.85%  | 0.00%    | 0.10%   |
| Language Identification (VoxForge)               | LLM-C↑ | 95.50%        | 13.50%  | 22.00%     | 8.50%       | 84.50%          | 93.00%               | 18.00% | 6.00%    | 0.00%   |
| Age Classification (Common Voice)                | LLM-C↑ | 0.50%         | 19.75%  | 22.75%     | 21.00%      | 26.00%          | 35.50%               | 21.00% | 23.50%   | 17.75%  |
| Gender Recognition by Voice (Common Voice)       | LLM-C↑ | 1.00%         | 81.50%  | 63.00%     | 57.00%      | 53.50%          | 97.50%               | 35.00% | 51.00%   | 63.00%  |
| HEAR Speaker Count Identification (LibriCount)   | NAR↓   | 30.40%        | 42.10%  | 35.90%     | 42.50%      | 3.50%           | 0.40%                | 22.00% | 23.20%   | 26.20%  |
| HEAR Speaker Count Identification (LibriCount)   | ACC↑   | 17.53%        | 16.93%  | 13.26%     | 17.22%      | 10.26%          | 13.76%               | 5.51%  | 13.80%   | 17.48%  |

Continued on next page

| Task  | Metric | Whisper-LLaMA | LTU-AS | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM | MU-LLaMA | GAMA-IT |
|---|--------|---------------|--------|------------|-------------|-----------------|----------------------|--------|----------|---------|
| Multi-Speaker Detection (VCTK)                    | LLM-C↑ | 46.50%        | 23.50% | 36.00%     | 30.50%      | 57.00%          | 51.00%               | 32.50% | 39.00%   | 48.50%  |
| Speaker Counting (LibriTTS-TestClean)             | LLM-C↑ | 18.50%        | 6.00%  | 14.50%     | 13.50%      | 17.50%          | 26.50%               | 14.00% | 4.00%    | 12.00%  |
| Speaker Verification (LibriSpeech-TestClean)      | LLM-C↑ | 48.00%        | 37.50% | 45.00%     | 54.50%      | 43.50%          | 52.00%               | 19.00% | 24.00%   | 13.00%  |
| Speaker Verification (LibriSpeech-TestOther)      | LLM-C↑ | 41.00%        | 36.00% | 51.00%     | 43.00%      | 49.00%          | 50.00%               | 20.00% | 27.00%   | 23.00%  |
| Speaker Verification (VCTK)                       | LLM-C↑ | 48.00%        | 38.50% | 54.00%     | 58.50%      | 32.50%          | 51.00%               | 24.00% | 19.50%   | 7.00%   |
| SuperbSD (Libri2Mix-Test)                         | NAR↓   | 93.00%        | 99.50% | 100.00%    | 100.00%     | 99.50%          | 98.00%               | 97.50% | 100.00%  | 60.50%  |
| SuperbSD (Libri2Mix-Test)                         | DER↓   | 74.81%        | 74.62% | -          | -           | 44.26%          | 93.29%               | 90.53% | -        | 74.12%  |
| SuperbSV (SuperbHiddenSet)                        | LLM-C↑ | 48.00%        | 47.00% | 45.00%     | 49.00%      | 51.00%          | 50.00%               | 55.00% | 44.00%   | 45.00%  |
| Noise Detection (LJSpeech, MUSAN-Gaussian)        | LLM-C↑ | 45.50%        | 18.50% | 50.00%     | 47.50%      | 42.00%          | 49.00%               | 39.00% | 41.00%   | 21.50%  |
| Noise Detection (LJSpeech, MUSAN-Music)           | LLM-C↑ | 52.00%        | 8.50%  | 52.00%     | 47.00%      | 50.50%          | 47.00%               | 44.50% | 33.00%   | 6.00%   |
| Noise Detection (LJSpeech, MUSAN-Noise)           | LLM-C↑ | 47.00%        | 13.50% | 50.50%     | 49.00%      | 50.50%          | 49.50%               | 44.50% | 48.00%   | 9.00%   |
| Noise Detection (LJSpeech, MUSAN-Speech)          | LLM-C↑ | 48.00%        | 13.50% | 53.00%     | 36.00%      | 50.50%          | 47.00%               | 45.00% | 50.50%   | 51.00%  |
| Noise Detection (VCTK, MUSAN-Music)               | LLM-C↑ | 46.00%        | 32.50% | 45.50%     | 57.50%      | 42.00%          | 54.50%               | 52.00% | 37.50%   | 3.50%   |
| Noise Detection (VCTK, MUSAN-Noise)               | LLM-C↑ | 44.50%        | 26.50% | 45.50%     | 55.00%      | 44.50%          | 61.00%               | 47.00% | 54.00%   | 4.50%   |
| Noise Detection (VCTK, MUSAN-Speech)              | LLM-C↑ | 45.00%        | 51.50% | 43.50%     | 52.50%      | 46.50%          | 57.00%               | 44.50% | 48.50%   | 56.50%  |
| Noise Detection (VCTK, MUSAN-Gaussian)            | LLM-C↑ | 53.00%        | 15.50% | 49.00%     | 56.50%      | 45.50%          | 54.50%               | 46.00% | 38.50%   | 11.50%  |
| Noise SNR Level Prediction (VCTK, MUSAN-Gaussian) | LLM-C↑ | 24.00%        | 8.00%  | 13.50%     | 15.50%      | 13.00%          | 17.50%               | 17.50% | 15.00%   | 1.50%   |
| Noise SNR Level Prediction (VCTK, MUSAN-Music)    | LLM-C↑ | 23.00%        | 9.00%  | 10.00%     | 8.50%       | 10.00%          | 14.00%               | 15.50% | 9.00%    | 1.50%   |
| Noise SNR Level Prediction (VCTK, MUSAN-Noise)    | LLM-C↑ | 23.00%        | 12.50% | 9.50%      | 16.00%      | 7.50%           | 19.00%               | 14.00% | 11.00%   | 0.50%   |

*Continued on next page*

| Task  | Metric | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM  | MU-LLaMA | GAMA-IT |
|---|--------|---------------|---------|------------|-------------|-----------------|----------------------|---------|----------|---------|
| Noise SNR Level Prediction (VCTK, MUSAN-Speech)           | LLM-C↑ | 23.50%        | 11.50%  | 13.00%     | 1.50%       | 14.00%          | 20.50%               | 15.00%  | 10.50%   | 0.50%   |
| Reverberation Detection (LJSpeech, RirsNoises-LargeRoom)  | LLM-C↑ | 44.50%        | 11.00%  | 19.00%     | 34.00%      | 28.00%          | 48.00%               | 42.50%  | 17.50%   | 6.00%   |
| Reverberation Detection (LJSpeech, RirsNoises-MediumRoom) | LLM-C↑ | 40.00%        | 4.50%   | 14.50%     | 25.00%      | 31.00%          | 48.00%               | 41.00%  | 10.50%   | 5.00%   |
| Reverberation Detection (LJSpeech, RirsNoises-SmallRoom)  | LLM-C↑ | 45.00%        | 4.50%   | 9.00%      | 19.50%      | 24.00%          | 48.00%               | 42.50%  | 14.00%   | 6.50%   |
| Reverberation Detection (VCTK, RirsNoises-LargeRoom)      | LLM-C↑ | 43.00%        | 24.50%  | 18.50%     | 25.00%      | 18.00%          | 46.00%               | 37.00%  | 14.50%   | 7.50%   |
| Reverberation Detection (VCTK, RirsNoises-MediumRoom)     | LLM-C↑ | 46.50%        | 22.50%  | 8.50%      | 19.50%      | 20.00%          | 46.00%               | 38.50%  | 14.00%   | 6.50%   |
| Reverberation Detection (VCTK, RirsNoises-SmallRoom)      | LLM-C↑ | 47.00%        | 18.50%  | 8.00%      | 13.00%      | 17.50%          | 46.00%               | 36.00%  | 10.50%   | 11.00%  |
| N-Best Correction (Librispeech-TestOther)                 | LLM-C↑ | 31.80%        | 22.80%  | 23.00%     | 29.20%      | 32.80%          | 29.00%               | 32.60%  | 20.80%   | 30.80%  |
| AAVE Speech Recognition (CORAAL)                          | WER↓   | 21.73%        | 97.34%  | 23.99%     | 31.56%      | 96.81%          | 38.55%               | 34.91%  | 136.50%  | 102.92% |
| Code-switch Speech Recognition (NTUML2021)                | MER↓   | 424.58%       | 215.51% | 293.50%    | 185.27%     | 165.54%         | 116.88%              | 172.03% | 130.18%  | 193.47% |
| Code-Switching Speech Recognition (ASCEND)                | NAR↓   | 68.20%        | 12.20%  | 5.40%      | 88.80%      | 1.60%           | 0.00%                | 30.40%  | 69.40%   | 14.60%  |
| Code-Switching Speech Recognition (ASCEND)                | ACC↑   | 28.30%        | 50.57%  | 58.14%     | 14.29%      | 44.31%          | 60.80%               | 10.92%  | 13.73%   | 9.84%   |
| Multi-Lingual Speech Recognition (MLS-de)                 | WER↓   | 73.80%        | 132.35% | 34.46%     | 25.21%      | 79.60%          | 24.96%               | 49.56%  | 99.37%   | 105.81% |
| Multi-Lingual Speech Recognition (MLS-en)                 | WER↓   | 65.03%        | 100.29% | 9.37%      | 9.47%       | 27.97%          | 17.47%               | 17.23%  | 96.03%   | 97.09%  |
| Multi-Lingual Speech Recognition (MLS-es)                 | WER↓   | 72.49%        | 123.26% | 23.45%     | 16.10%      | 58.82%          | 18.19%               | 44.90%  | 103.19%  | 101.75% |
| Multi-Lingual Speech Recognition (MLS-fr)                 | WER↓   | 71.09%        | 133.86% | 26.52%     | 21.27%      | 46.27%          | 19.01%               | 42.35%  | 102.25%  | 102.01% |

*Continued on next page*

| Task  | Metric | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM  | MU-LLaMA | GAMA-IT  |
|---|--------|---------------|---------|------------|-------------|-----------------|----------------------|---------|----------|----------|
| Multi-Lingual Speech Recognition (MLS-it)       | WER↓   | 85.38%        | 125.39% | 39.11%     | 31.54%      | 56.55%          | 33.16%               | 46.14%  | 102.72%  | 108.84%  |
| Multi-Lingual Speech Recognition (MLS-nl)       | WER↓   | 74.96%        | 120.29% | 32.11%     | 29.37%      | 140.01%         | 42.94%               | 55.56%  | 102.05%  | 101.22%  |
| Multi-Lingual Speech Recognition (MLS-pl)       | WER↓   | 76.79%        | 139.04% | 73.83%     | 51.53%      | 278.19%         | 101.41%              | 90.11%  | 100.23%  | 105.43%  |
| Multi-Lingual Speech Recognition (MLS-pt)       | WER↓   | 73.63%        | 136.66% | 31.88%     | 25.07%      | 117.33%         | 24.72%               | 46.26%  | 99.71%   | 103.48%  |
| PTBR Speech Recognition (Common Voice)          | WER↓   | 45.75%        | 248.05% | 69.88%     | 73.91%      | 233.85%         | 38.57%               | 62.93%  | 639.12%  | 270.64%  |
| SuperbASR (LibriSpeech-TestClean)               | WER↓   | 33.96%        | 96.28%  | 15.11%     | 2.79%       | 69.35%          | 36.70%               | 6.87%   | 103.67%  | 116.66%  |
| SuperbASR (LibriSpeech-TestOther)               | WER↓   | 42.14%        | 91.89%  | 11.44%     | 4.31%       | 79.53%          | 40.28%               | 9.17%   | 111.03%  | 130.00%  |
| SuperbOOD Asr Ar (Common Voice)                 | WER↓   | 51.04%        | 245.85% | 216.51%    | 178.02%     | 289.25%         | 178.21%              | 149.15% | 225.28%  | 504.53%  |
| SuperbOOD Asr Es (Common Voice)                 | WER↓   | 10.93%        | 150.00% | 98.96%     | 99.22%      | 100.83%         | 75.08%               | 99.38%  | 141.68%  | 303.28%  |
| SuperbOOD Asr Spon (CHIME6-Test)                | WER↓   | 65.71%        | 122.86% | 61.11%     | 62.80%      | 92.17%          | 80.21%               | 136.24% | 137.25%  | 266.72%  |
| SuperbOOD Asr Zh (Common Voice)                 | CER↓   | 29.67%        | 508.02% | 609.28%    | 310.33%     | 449.32%         | 270.63%              | 445.36% | 435.33%  | 1097.41% |
| Target Speaker ASR (AMI-test)                   | WER↓   | 143.02%       | 133.86% | 273.90%    | 187.15%     | 207.01%         | 132.13%              | 266.26% | 108.32%  | 140.37%  |
| Multi-Speaker Detection (LibriSpeech-TestClean) | LLM-C↑ | 46.00%        | 29.00%  | 16.50%     | 22.00%      | 58.00%          | 52.00%               | 22.00%  | 42.50%   | 53.00%   |
| Speech Command Recognition (AudioMNIST)         | NAR↓   | 4.80%         | 76.40%  | 3.07%      | 0.40%       | 0.67%           | 0.00%                | 0.67%   | 1.60%    | 76.40%   |
| Speech Command Recognition (AudioMNIST)         | ACC↑   | 88.80%        | 77.40%  | 96.70%     | 74.43%      | 96.24%          | 77.07%               | 93.42%  | 9.89%    | 8.47%    |
| Speech Text Matching (LibriSpeech-TestClean)    | LLM-C↑ | 86.50%        | 43.50%  | 57.00%     | 59.50%      | 64.00%          | 92.00%               | 52.00%  | 44.50%   | 42.00%   |
| Speech Text Matching (LibriSpeech-TestOther)    | LLM-C↑ | 80.50%        | 42.50%  | 56.50%     | 60.00%      | 66.00%          | 92.50%               | 53.00%  | 47.00%   | 40.00%   |
| Speech Text Matching (LJSpeech)                 | LLM-C↑ | 83.50%        | 44.50%  | 57.00%     | 60.00%      | 67.00%          | 90.00%               | 52.00%  | 37.50%   | 41.00%   |

Continued on next page

| Task  | Metric | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM | MU-LLaMA | GAMA-IT |
|---|--------|---------------|---------|------------|-------------|-----------------|----------------------|--------|----------|---------|
| Spoken Term Detection (LibriSpeech-TestClean) | LLM-C↑ | 76.50%        | 28.00%  | 60.00%     | 54.50%      | 77.00%          | 61.50%               | 51.50% | 24.00%   | 37.50%  |
| Spoken Term Detection (LibriSpeech-TestOther) | LLM-C↑ | 75.50%        | 23.00%  | 54.50%     | 50.50%      | 72.50%          | 64.00%               | 46.50% | 25.50%   | 34.00%  |
| Spoken Term Detection (LJSpeech)              | LLM-C↑ | 83.50%        | 33.50%  | 57.00%     | 53.50%      | 74.50%          | 79.50%               | 51.00% | 25.00%   | 25.00%  |
| SuperbKS (Speech Commands V1-Test)            | LLM-C↑ | 36.50%        | 1.00%   | 30.50%     | 2.00%       | 60.50%          | 47.00%               | 43.00% | 4.00%    | 2.00%   |
| SuperbQbE (Quesst14-Eval)                     | LLM-C↑ | 46.50%        | 45.00%  | 49.00%     | 51.50%      | 48.00%          | 53.50%               | 49.50% | 51.00%   | 2.50%   |
| Stuttering Detection (SEP28k)                 | LLM-C↑ | 49.10%        | 51.00%  | 50.50%     | 50.30%      | 50.50%          | 52.40%               | 55.00% | 49.60%   | 47.40%  |
| Voice Disorder Classification (VOICED)        | LLM-C↑ | 13.46%        | 1.92%   | 13.46%     | 17.31%      | 13.46%          | 16.35%               | 18.27% | 21.15%   | 6.73%   |
| Conversation Matching                         | LLM-C↑ | 77.78%        | 5.56%   | 57.41%     | 37.04%      | 51.85%          | 66.67%               | 62.96% | 3.70%    | 24.07%  |
| Dialogue Act Classification (SLUE-HVB)        | LLM-C↑ | 11.08%        | 0.00%   | 36.00%     | 0.92%       | 5.00%           | 15.83%               | 13.17% | 1.83%    | 1.33%   |
| Dialogue Act Classification (Daily Talk)      | LLM-C↑ | 29.00%        | 10.50%  | 34.00%     | 40.50%      | 42.00%          | 30.00%               | 36.50% | 18.50%   | 4.00%   |
| Dialogue Act Pairing (Daily Talk)             | LLM-C↑ | 50.00%        | 3.50%   | 51.00%     | 48.00%      | 43.50%          | 46.00%               | 37.00% | 40.00%   | 25.00%  |
| SuperbIC (SLURP)                              | LLM-C↑ | 36.50%        | 0.06%   | 21.00%     | 5.50%       | 26.50%          | 26.00%               | 28.00% | 3.50%    | 1.00%   |
| Intent Classification (SLURP-Action)          | LLM-C↑ | 27.00%        | 14.00%  | 28.50%     | 28.50%      | 44.00%          | 68.50%               | 37.50% | 6.50%    | 1.50%   |
| Intent Classification (SLURP-Intent)          | LLM-C↑ | 45.50%        | 11.50%  | 27.00%     | 19.00%      | 68.50%          | 56.00%               | 56.50% | 4.00%    | 0.50%   |
| Named Entity Localization (SLUE-VoxPopuli)    | NAR↓   | 81.25%        | 100.00% | 100.00%    | 100.00%     | 100.00%         | 25.00%               | 93.75% | 100.00%  | 100.00% |
| Named Entity Localization (SLUE-VoxPopuli)    | F1↑    | 1.4006        | -       | -          | -           | -               | 1.1063               | 0.8681 | -        | -       |
| Named Entity Recognition (SLUE-VoxPopuli)     | NAR↓   | 90.76%        | 98.37%  | 95.65%     | 98.37%      | 95.11%          | 91.30%               | 91.85% | 97.83%   | 96.74%  |
| Named Entity Recognition (SLUE-VoxPopuli)     | IoU↑   | 0             | 0.3333  | 0          | 0           | 0               | 0.0625               | 0.0222 | 0        | 0       |
| Sarcasm Detection (Mustard)                   | LLM-C↑ | 44.50%        | 15.00%  | 38.50%     | 46.00%      | 42.50%          | 48.00%               | 46.50% | 43.00%   | 8.00%   |

Continued on next page

| Task  | Metric          | Whisper-LLaMA | LTU-AS | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM  | MU-LLaMA | GAMA-IT |
|---|-----------------|---------------|--------|------------|-------------|-----------------|----------------------|---------|----------|---------|
| Semantic Textual Similarity (SpokenSTS)                               | LLM-C↑          | 48.80%        | 47.20% | 47.20%     | 50.40%      | 42.80%          | 46.80%               | 50.40%  | 51.20%   | 38.80%  |
| Speech Sentiment Analysis (MELD)                                      | LLM-C↑          | 38.50%        | 37.50% | 33.50%     | 35.09%      | 38.50%          | 49.00%               | 44.00%  | 27.50%   | 27.00%  |
| Spoken Digit Arithmetic (AudioMNIST)                                  | LLM-C↑          | 43.50%        | 13.00% | 31.00%     | 33.00%      | 15.00%          | 43.50%               | 40.50%  | 4.50%    | 13.00%  |
| SuperbSF (AudioSnips-Test)  | NAR↓            | 38.00%        | 88.50% | 34.00%     | 68.50%      | 51.00%          | 45.00%               | 44.00%  | 99.00%   | 97.00%  |
| SuperbSF (AudioSnips-Test)  | Slot Type FI↑   | 0.8637        | 0.8864 | 0.9325     | 0.9396      | 0.8616          | 0.7704               | 0.907   | 1        | 0.8939  |
| SuperbSF (AudioSnips-Test)  | Slot Value CER↓ | 0.3178        | 0.6432 | 0.4059     | 0.3349      | 0.4417          | 0.3065               | 0.4752  | 1.587    | 1.425   |
| Code-Switching Semantic Grammar Acceptability Comparison (CSZS-zh-en) | LLM-C↑          | 51.50%        | 45.00% | 50.00%     | 49.50%      | 27.50%          | 49.50%               | 27.50%  | 28.00%   | 19.00%  |
| Nonce Word Detection (sWUGGY)   | LLM-C↑          | 48.43%        | 30.20% | 48.43%     | 48.43%      | 21.08%          | 50.14%               | 48.15%  | 35.04%   | 19.37%  |
| PoS Estimation (LibriTTS)   | POS↓            | 2.5722        | 1.2792 | 2.9675     | 1.199       | 3.062           | 2.3091               | 3.1126  | 1.5327   | 1.9635  |
| PoS Estimation with transcription (LibriTTS)                          | POS↓            | 0.9081        | 1.3507 | 1.7583     | 1.0601      | 2.158           | 1.7514               | 2.0463  | 1.2394   | 1.3911  |
| Sentence Grammar Acceptability (sBLIMP)                               | LLM-C↑          | 52.78%        | 6.75%  | 27.38%     | 42.06%      | 28.97%          | 49.40%               | 49.80%  | 50.79%   | 3.57%   |
| SuperbST (CoVoST2-Test)   | Sacre Bleu↑     | 17.372        | 0.107  | 18.7884    | 16.7835     | 6.9846          | 23.3458              | 21.7423 | 0.0566   | 0.0219  |

Table 6: Task-level evaluation results using LLaMA-3.1-8B-Instruct.

| Task                                    | Metric | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM | MU-LLaMA | GAMA-IT |
|---|--------|---------------|---------|------------|-------------|-----------------|----------------------|--------|----------|---------|
| Chord Classification                    | LLM-C↑ | 4.00%         | 5.00%   | 7.50%      | 13.00%      | 43.00%          | 16.00%               | 45.50% | 50.00%   | 0.50%   |
| MARBLE Key Detection (GiantSteps - Key) | LLM-C↑ | 3.50%         | 1.00%   | 1.50%      | 1.00%       | 3.00%           | 13.00%               | 2.00%  | 0.50%    | 11.50%  |
| HEAR Music Transcription (MAESTRO)      | NAR↓   | 95.14%        | 100.00% | 100.00%    | 100.00%     | 97.30%          | 96.22%               | 84.86% | 100.00%  | 100.00% |
| HEAR Music Transcription (MAESTRO)      | TER↓   | 99.45%        | -       | -          | -           | 99.54%          | 100.00%              | 99.48% | -        | -       |

*Continued on next page*

| Task  | Metric | Whisper-LLaMA | LFU-AS | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM | MU-LLaMA | GAMA-IT |
|---|--------|---------------|--------|------------|-------------|-----------------|----------------------|--------|----------|---------|
| HEAR Percussion Instruments Tonic Classification (Mridangam Tonic)    | LLM-C↑ | 2.40%         | 3.30%  | 14.00%     | 15.10%      | 14.00%          | 13.20%               | 16.00% | 14.30%   | 8.60%   |
| Instrument Pitch Classification (Nsynth)                              | LLM-C↑ | 0.70%         | 0.90%  | 1.70%      | 5.40%       | 1.60%           | 18.90%               | 1.20%  | 0.40%    | 1.00%   |
| Pitch Extraction By Lyrics (CSD)                                      | NAR↓   | 30.00%        | 73.50% | 57.00%     | 85.50%      | 76.50%          | 73.00%               | 57.50% | 97.50%   | 53.00%  |
| Pitch Extraction By Lyrics (CSD)                                      | TER↓   | 1.86          | 2.82   | 4.32       | 7.38        | 2.33            | 2.42                 | 2.71   | 1.00     | 3.87    |
| Emotion Classification In Songs (EMOTIFY)                             | LLM-C↑ | 10.00%        | 0.00%  | 1.70%      | 3.30%       | 21.70%          | 8.30%                | 8.30%  | 1.70%    | 1.70%   |
| MARBLE Emotion Detection (MTG)  | LLM-C↑ | 3.60%         | 8.70%  | 0.80%      | 3.70%       | 9.90%           | 3.30%                | 1.80%  | 7.20%    | 1.10%   |
| HEAR Music Genre Classification (ISMIR04)                             | LLM-C↑ | 13.07%        | 10.05% | 26.13%     | 45.73%      | 31.66%          | 37.19%               | 15.58% | 26.13%   | 1.51%   |
| MARBLE Genre Classification (MTG)                                     | LLM-C↑ | 1.20%         | 2.80%  | 0.90%      | 1.50%       | 5.60%           | 0.50%                | 1.20%  | 1.80%    | 0.40%   |
| MARBLE Music Tagging (MagnaTagATune)                                  | LLM-C↑ | 7.50%         | 3.00%  | 7.00%      | 16.00%      | 5.00%           | 3.50%                | 14.50% | 9.00%    | 7.00%   |
| MARBLE Music Tagging (MTG)  | LLM-C↑ | 5.20%         | 3.00%  | 1.40%      | 0.90%       | 4.50%           | 4.80%                | 2.60%  | 3.80%    | 1.10%   |
| Music Genre Classification (FMA)                                      | LLM-C↑ | 13.40%        | 1.80%  | 9.80%      | 16.10%      | 8.00%           | 17.00%               | 8.90%  | 9.80%    | 5.40%   |
| HEAR Percussion Instruments Classification (Beijing Opera Percussion) | LLM-C↑ | 2.98%         | 3.82%  | 10.60%     | 20.74%      | 16.48%          | 29.26%               | 25.42% | 22.02%   | 11.84%  |
| HEAR Percussion Instruments Stroke Classification (Mridangam Stroke)  | LLM-C↑ | 0.10%         | 5.00%  | 7.90%      | 4.90%       | 12.00%          | 11.50%               | 14.40% | 10.60%   | 0.90%   |
| Instrument Classification (Nsynth)                                    | LLM-C↑ | 3.70%         | 12.90% | 14.80%     | 15.30%      | 21.30%          | 55.40%               | 13.20% | 10.90%   | 9.60%   |
| Instrument Combination Recognition (OpenMIC-2018)                     | LLM-C↑ | 8.10%         | 12.00% | 21.90%     | 27.90%      | 33.60%          | 25.50%               | 24.90% | 22.80%   | 10.20%  |
| Instrument Source Classification (Nsynth)                             | LLM-C↑ | 9.90%         | 36.20% | 39.20%     | 48.00%      | 37.10%          | 37.90%               | 38.00% | 41.40%   | 47.70%  |
| MARBLE Instrument Classification (MTG)                                | LLM-C↑ | 1.50%         | 6.90%  | 0.20%      | 0.00%       | 6.70%           | 2.60%                | 1.30%  | 3.00%    | 7.60%   |
| Singing Automatic MOS Prediction (SingMOS)                            | NAR↓   | 99.50%        | 68.22% | 0.33%      | 16.14%      | 17.30%          | 0.83%                | 23.13% | 97.84%   | 98.50%  |

Continued on next page



| Task  | Metric     | Whisper-LLaMA | LTU-AS    | SALMONN 7B  | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM     | MU-LLaMA   | GAMA-IT  |
|---|------------|---------------|-----------|-------------|-------------|-----------------|----------------------|------------|------------|----------|
| Singing Automatic MOS Prediction (SingMOS)                      | MSE↓       | 6.5670        | 2.4110    | 5.3810      | 33.1090     | 14.8900         | 1.3810               | 11.2760    | 4.3720     | 2.3070   |
| Singing Automatic MOS Prediction (SingMOS)                      | KTAU↑      | 0.3330        | 0.0090    | 0.2330      | 0.0230      | 0.0710          | 0.0180               | 0.0300     | 0.1640     | 0.2250   |
| Singing Automatic MOS Prediction (SingMOS)                      | LCC↑       | 0.6150        | -0.0780   | 0.0240      | 0.0140      | 0.0820          | 0.0400               | -0.0050    | 0.0980     | 0.2330   |
| Singing Automatic MOS Prediction (SingMOS)                      | SRCC↑      | 0.5000        | 0.0120    | 0.2210      | 0.0330      | 0.0870          | 0.0230               | 0.0380     | 0.1870     | 0.2790   |
| MARBLE Beat Tracking (ASAP)                                     | NAR↓       | 100.00%       | 100.00%   | 100.00%     | 100.00%     | 100.00%         | 100.00%              | 100.00%    | 100.00%    | 100.00%  |
| MARBLE Beat Tracking (ASAP)                                     | ERR↓       | -             | -         | -           | -           | -               | -                    | -          | -          | -        |
| Music Beat Tracking (ASAP)                                      | NAR↓       | 53.64%        | 91.82%    | 50.00%      | 65.45%      | 51.82%          | 40.91%               | 43.64%     | 78.18%     | 95.45%   |
| Music Beat Tracking (ASAP)                                      | Miss Time↓ | 59.7630       | 63.6760   | 98.6640     | 56.3060     | 131.7640        | 31.3990              | 81.3490    | 199.8770   | 22.2000  |
| Audio Editing Identification (People's Speech)                  | NAR↓       | 87.50%        | 93.75%    | 84.38%      | 84.38%      | 46.88%          | 31.25%               | 81.25%     | 87.50%     | 37.50%   |
| Audio Editing Identification (People's Speech)                  | ACC↑       | 153.6830      | -650.0000 | 283.7050    | 82.8400     | 123.3700        | 165.0680             | -49.7720   | 154.1090   | 147.6570 |
| Scene Fake Detection (ASPIRE)                                   | LLM-C↑     | 29.00%        | 48.00%    | 48.00%      | 52.00%      | 49.00%          | 54.00%               | 49.00%     | 19.00%     | 47.00%   |
| Audio Deep Fake Detection (LJSpeech, WaveFake, MUSDB18HQ)       | LLM-C↑     | 18.50%        | 34.20%    | 29.10%      | 42.00%      | 36.00%          | 39.70%               | 42.00%     | 40.30%     | 37.90%   |
| Singing Voice Deepfake Detection (CtrSVDD, ACEKiSing, M4Singer) | LLM-C↑     | 10.00%        | 35.20%    | 29.80%      | 38.60%      | 29.70%          | 31.90%               | 27.10%     | 40.80%     | 22.20%   |
| Audio Duration Prediction (NTUML2021)                           | NAR↓       | 33.00%        | 28.00%    | 0.50%       | 49.00%      | 0.00%           | 0.00%                | 0.00%      | 46.00%     | 14.00%   |
| Audio Duration Prediction (NTUML2021)                           | MSE↓       | 13.3660       | 35.6460   | 98,768.7290 | 2,005.3330  | 28.9400         | 62.1550              | 3,243.7400 | 1,779.0420 | 62.7500  |
| HEAR Music Speech Classification (MAESTRO, Librispeech)         | LLM-C↑     | 65.00%        | 68.50%    | 87.50%      | 99.50%      | 94.00%          | 89.50%               | 22.00%     | 53.50%     | 88.00%   |
| Sound Effect Detection (RemFx)                                  | LLM-C↑     | 4.00%         | 7.67%     | 14.67%      | 18.00%      | 17.00%          | 19.67%               | 17.33%     | 17.00%     | 16.33%   |
| Speech Detection (LibriSpeech-TestClean)                        | LLM-C↑     | 48.00%        | 32.00%    | 46.00%      | 66.00%      | 45.50%          | 56.50%               | 36.50%     | 46.50%     | 56.00%   |
| Speech Detection (LibriSpeech-TestOther)                        | LLM-C↑     | 51.00%        | 33.00%    | 46.00%      | 66.50%      | 45.00%          | 55.50%               | 38.50%     | 45.00%     | 51.00%   |

*Continued on next page*

| Task   | Metric          | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM  | MU-LLaMA | GAMA-IT |
|--|-----------------|---------------|---------|------------|-------------|-----------------|----------------------|---------|----------|---------|
| Speech Detection (LJSpeech)                                      | LLM-C↑          | 55.00%        | 46.00%  | 60.50%     | 74.00%      | 56.50%          | 47.50%               | 51.50%  | 49.50%   | 76.00%  |
| Children Song Transcript Verification (CSD)                      | WER↓            | 58.70%        | 102.70% | 104.00%    | 101.00%     | 128.70%         | 94.90%               | 155.50% | 100.00%  | 100.10% |
| Lyric Translation (SingSet)                                      | Sacre Bleu↑     | 2.1090        | 1.0360  | 2.8130     | 0.9940      | 3.3970          | 3.9160               | 1.8970  | 0.1590   | 0.0740  |
| Song Lyric Recognition (SingSet)                                 | MER↓            | 1.0193        | 1.6453  | 3.5636     | 2.6600      | 1.1841          | 1.1562               | 4.8113  | 1.1753   | 1.6852  |
| MARBLE Vocal Technique Detection (VocalSet)                      | LLM-C↑          | 8.50%         | 2.00%   | 12.00%     | 11.00%      | 7.50%           | 15.50%               | 10.00%  | 9.00%    | 4.00%   |
| Audio Segment Retrieval (Clotho)                                 | NAR↓            | 86.00%        | 92.00%  | 24.00%     | 24.00%      | 42.00%          | 28.00%               | 44.00%  | 84.00%   | 42.00%  |
| Audio Segment Retrieval (Clotho)                                 | IoU↑            | 0.3700        | 0.2970  | 0.1900     | 0.2460      | 0.0570          | 0.1930               | 0.1710  | 0.0000   | 0.1590  |
| HEAR Sound Event Detection (DCASE2016 Task2)                     | LLM-C↑          | 0.00%         | 0.00%   | 28.60%     | 35.70%      | 0.00%           | 7.10%                | 42.90%  | 7.10%    | 14.30%  |
| Multi-channel Sound Event Understanding (STARSS23)               | LLM-C↑          | 36.60%        | 9.20%   | 36.60%     | 26.10%      | 38.00%          | 43.00%               | 23.90%  | 8.50%    | 15.50%  |
| Animal Classification (WaveSource-Test)                          | LLM-C↑          | 9.80%         | 36.30%  | 58.50%     | 66.50%      | 75.80%          | 34.50%               | 13.30%  | 4.30%    | 15.30%  |
| Bird Sound Detection (Warblr10k)                                 | LLM-C↑          | 32.50%        | 41.50%  | 75.50%     | 77.50%      | 78.50%          | 82.50%               | 75.00%  | 45.00%   | 35.50%  |
| Cat Emotion Classification (Cat Sound Classification Dataset V2) | LLM-C↑          | 10.00%        | 2.00%   | 2.00%      | 4.00%       | 4.00%           | 16.00%               | 8.00%   | 2.00%    | 8.00%   |
| Cornell Birdcall Identification                                  | LLM-C↑          | 10.00%        | 3.30%   | 0.00%      | 3.30%       | 10.00%          | 10.00%               | 13.30%  | 3.30%    | 13.30%  |
| Environmental Sound Classification (ESC50-Animals)               | LLM-C↑          | 14.00%        | 17.50%  | 20.00%     | 8.00%       | 78.00%          | 34.00%               | 5.50%   | 2.50%    | 7.50%   |
| HEAR Beehive States Classification (Beehive States)              | LLM-C↑          | 92.50%        | 70.00%  | 60.50%     | 92.00%      | 56.00%          | 92.00%               | 82.50%  | 61.00%   | 70.00%  |
| Domestic Environment Sound Event Detection (DESED-PublicEval)    | NAR↓            | 97.22%        | 100.00% | 99.72%     | 100.00%     | 99.44%          | 99.72%               | 100.00% | 98.33%   | 99.72%  |
| Domestic Environment Sound Event Detection (DESED-PublicEval)    | Event-based F1↑ | 0.4800        | 0.0000  | 0.0000     | 0.0000      | 0.1680          | 1.0000               | 0.0000  | 1.0000   | 0.0000  |

Continued on next page

| Task  | Metric      | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM  | MU-LLaMA | GAMA-IT |
|---|-------------|---------------|---------|------------|-------------|-----------------|----------------------|---------|----------|---------|
| Emergency Traffic Detection (Large-Scale-Audio-dataset)                         | LLM-C↑      | 53.50%        | 32.50%  | 57.50%     | 85.00%      | 66.00%          | 75.50%               | 40.50%  | 46.00%   | 73.00%  |
| Environmental Sound Classification (ESC50-Exterior And Urban Noises)            | LLM-C↑      | 7.00%         | 11.00%  | 4.50%      | 5.00%       | 67.50%          | 25.00%               | 6.50%   | 1.00%    | 12.00%  |
| Environmental Sound Classification (ESC50-Interior And Domestic Sounds)         | LLM-C↑      | 13.50%        | 7.50%   | 1.00%      | 0.50%       | 56.50%          | 15.00%               | 10.00%  | 0.00%    | 6.50%   |
| Environmental Sound Classification (ESC50-Natural Soundscapes And Water Sounds) | LLM-C↑      | 5.00%         | 10.00%  | 7.00%      | 1.00%       | 68.50%          | 8.00%                | 8.50%   | 0.00%    | 9.00%   |
| Environmental Sound Classification (UrbanSound8K-Urban Noises)                  | LLM-C↑      | 13.00%        | 5.80%   | 2.30%      | 5.30%       | 35.10%          | 7.90%                | 10.90%  | 0.50%    | 4.70%   |
| Environment Recognition (ESC50)   | LLM-C↑      | 7.00%         | 45.60%  | 40.40%     | 21.10%      | 68.40%          | 66.70%               | 38.60%  | 10.50%   | 52.60%  |
| HEAR Environmental Sound Classification (ESC50)                                 | LLM-C↑      | 10.70%        | 31.40%  | 64.00%     | 76.60%      | 73.00%          | 52.00%               | 22.80%  | 4.30%    | 43.40%  |
| HEAR Vocal Imitation Classification (Vocal Imitations)                          | LLM-C↑      | 15.33%        | 9.00%   | 2.00%      | 2.17%       | 17.83%          | 11.33%               | 22.83%  | 4.17%    | 7.67%   |
| Audio Spatial Distance Prediction (Spatial LibriSpeech)                         | NAR↓        | 45.00%        | 64.00%  | 81.00%     | 95.50%      | 80.00%          | 41.50%               | 39.50%  | 92.50%   | 55.50%  |
| Audio Spatial Distance Prediction (Spatial LibriSpeech)                         | MEDAE↓      | 0.7110        | 0.6800  | 0.6840     | 0.0000      | 1.3130          | 0.7840               | 0.7450  | 0.6340   | 0.5910  |
| How Far Are You (3DSpeaker)   | LLM-C↑      | 27.50%        | 13.00%  | 19.00%     | 8.00%       | 30.00%          | 26.50%               | 32.00%  | 19.00%   | 2.50%   |
| Sound Position Prediction   | NAR↓        | 100.00%       | 100.00% | 100.00%    | 100.00%     | 100.00%         | 100.00%              | 100.00% | 100.00%  | 100.00% |
| Sound Position Prediction   | Angle Diff↓ | -             | -       | -          | -           | -               | -                    | -       | -        | -       |
| Dialogue Emotion Classification (Daily Talk)                                    | LLM-C↑      | 54.00%        | 18.00%  | 28.50%     | 13.00%      | 33.50%          | 54.00%               | 22.50%  | 42.00%   | 41.00%  |
| Emoji Grounded Speech Emotion Recognition (RAVDESS)                             | LLM-C↑      | 1.10%         | 2.00%   | 2.20%      | 0.20%       | 3.50%           | 1.80%                | 0.10%   | 2.20%    | 5.00%   |
| Emotion Change Detection (Ravdess)  | LLM-C↑      | 5.00%         | 10.30%  | 10.00%     | 11.20%      | 1.00%           | 10.30%               | 3.90%   | 11.30%   | 2.90%   |

Continued on next page

| Task   | Metric    | Whisper-LLaMA | LTU-AS     | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM      | MU-LLaMA | GAMA-IT |
|--|-----------|---------------|------------|------------|-------------|-----------------|----------------------|-------------|----------|---------|
| Emotion Recognition (MELD)   | LLM-C↑    | 33.50%        | 14.00%     | 20.00%     | 10.00%      | 43.50%          | 36.50%               | 40.50%      | 1.50%    | 8.50%   |
| HEAR Emotion Recognition (CREMAD)  | LLM-C↑    | 10.60%        | 10.40%     | 18.80%     | 18.10%      | 64.00%          | 62.40%               | 26.30%      | 13.60%   | 16.00%  |
| SuperbER (RAVDESS)   | LLM-C↑    | 12.92%        | 24.58%     | 13.75%     | 13.33%      | 70.42%          | 76.25%               | 12.50%      | 13.75%   | 5.00%   |
| Covid19 Cough Audio Classification (CoughVid)                                    | LLM-C↑    | 40.00%        | 18.80%     | 4.90%      | 86.40%      | 9.30%           | 8.30%                | 18.60%      | 42.10%   | 6.00%   |
| Environmental Sound Classification (ESC50-Human And Non Speech Sounds)           | LLM-C↑    | 17.00%        | 11.00%     | 9.50%      | 3.00%       | 77.00%          | 29.00%               | 19.50%      | 0.00%    | 11.00%  |
| Human Non Speech Sound Recognition (Nonspeech7k, CommonVoice)                    | LLM-C↑    | 3.60%         | 5.70%      | 45.00%     | 30.70%      | 27.90%          | 23.60%               | 18.60%      | 17.10%   | 28.60%  |
| Human Screaming Detection (Environmentdb)  | LLM-C↑    | 52.50%        | 45.00%     | 52.50%     | 65.00%      | 62.50%          | 87.50%               | 47.50%      | 47.50%   | 67.50%  |
| Vocal Sound Recognition (VocalSound)   | LLM-C↑    | 32.50%        | 2.50%      | 11.70%     | 4.00%       | 75.70%          | 38.30%               | 22.10%      | 1.00%    | 3.90%   |
| Phoneme Segment Counting (Librispeech-words)                                     | NAR↓      | 23.29%        | 60.40%     | 37.69%     | 18.67%      | 9.52%           | 13.44%               | 26.81%      | 4.97%    | 89.74%  |
| Phoneme Segment Counting (Librispeech-words)                                     | ACC↑      | 0.1510        | 0.1400     | 0.1760     | 0.1190      | 0.0540          | 0.1530               | 0.1380      | 0.0070   | 0.1020  |
| Phoneme Segment Counting (Librispeech-words)                                     | Abs Diff↓ | 2.4090        | 79.7600    | 367.5620   | 38.1450     | 7.1110          | 2.9900               | 8.8670      | 15.2140  | 4.5020  |
| Phone Segment Counting (VoxAngeles)  | NAR↓      | 23.29%        | 70.96%     | 32.22%     | 14.07%      | 12.30%          | 28.23%               | 18.81%      | 46.88%   | 24.10%  |
| Phone Segment Counting (VoxAngeles)  | ACC↑      | 0.1510        | 0.1180     | 0.1910     | 0.2300      | 0.0730          | 0.0950               | 0.1610      | 0.1100   | 0.0430  |
| Phone Segment Counting (VoxAngeles)  | Abs Diff↓ | 2.2150        | 5,532.7270 | 8.1900     | 6.2270      | 2.8590          | 2.7220               | 51,811.8860 | 3.9190   | 2.8570  |
| SuperbPR (LibriSpeech-TestClean)   | PER↓      | 100.12%       | 102.75%    | 25.36%     | 24.60%      | 100.58%         | 100.96%              | 99.99%      | 110.27%  | 100.37% |
| SuperbPR (LibriSpeech-TestOther)   | PER↓      | 100.17%       | 101.08%    | 22.79%     | 22.91%      | 100.62%         | 100.88%              | 99.99%      | 111.18%  | 100.61% |
| Phonological Feature Classification (VoxAngeles-Consonant Place Of Articulation) | LLM-C↑    | 25.90%        | 2.20%      | 2.00%      | 5.10%       | 4.40%           | 9.90%                | 3.60%       | 3.10%    | 1.00%   |

Continued on next page

| Task  | Metric | Whisper-LLaMA | LTU-AS | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM | MU-LLaMA | GAMA-IT |
|---|--------|---------------|--------|------------|-------------|-----------------|----------------------|--------|----------|---------|
| Phonological Feature Classification (VoxAngeles-Manner Of Articulation) | LLM-C↑ | 16.90%        | 9.50%  | 2.20%      | 1.50%       | 5.50%           | 6.50%                | 7.00%  | 8.00%    | 6.00%   |
| Phonological Feature Classification (VoxAngeles-Phone)                  | LLM-C↑ | 5.90%         | 0.00%  | 2.80%      | -           | 9.50%           | 5.20%                | 79.90% | 0.60%    | 0.10%   |
| Phonological Feature Classification (VoxAngeles-Vowel Frontness)        | LLM-C↑ | 50.30%        | 14.90% | 41.30%     | 38.30%      | 50.50%          | 58.20%               | 47.00% | 42.20%   | 16.10%  |
| Phonological Feature Classification (VoxAngeles-Vowel Height)           | LLM-C↑ | 21.20%        | 15.10% | 32.00%     | 35.20%      | 30.80%          | 38.70%               | 39.30% | 39.10%   | 25.70%  |
| Phonological Feature Classification (VoxAngeles-Vowel Roundedness)      | LLM-C↑ | 41.10%        | 18.10% | 27.30%     | 61.30%      | 46.40%          | 70.10%               | 37.10% | 21.80%   | 17.90%  |
| Heteronym Differentiation (HeteronymEn)                                 | LLM-C↑ | 55.00%        | 25.00% | 42.00%     | 53.00%      | 45.00%          | 44.00%               | 49.00% | 29.00%   | 20.00%  |
| L2 English Accuracy (speechocean762-Ranking)                            | LLM-C↑ | 22.20%        | 30.80% | 50.00%     | 51.10%      | 35.30%          | 50.00%               | 37.20% | 51.90%   | 30.80%  |
| L2 English Accuracy (speechocean762-scoring)                            | NAR↓   | 15.44%        | 78.52% | 0.27%      | 32.48%      | 10.47%          | 0.94%                | 12.89% | 61.88%   | 15.84%  |
| L2 English Accuracy (speechocean762-scoring)                            | PCC↑   | 0.0260        | 0.2990 | 0.0670     | 0.0920      | 0.0250          | -0.0110              | 0.0930 | 0.0370   | -0.0310 |
| Multi-lingual Pronunciation Similarity (VoxAngeles)                     | LLM-C↑ | 37.90%        | 6.60%  | 46.60%     | 46.20%      | 14.80%          | 43.20%               | 44.40% | 20.20%   | 29.70%  |
| Accent Classification (Accentdb Extended)                               | LLM-C↑ | 17.50%        | 8.50%  | 6.00%      | 7.00%       | 24.00%          | 14.00%               | 7.50%  | 37.00%   | 4.00%   |
| Stress Detection (MIRSD)  | LLM-C↑ | 15.00%        | 6.50%  | 7.00%      | 18.50%      | 9.50%           | 20.00%               | 29.00% | 5.00%    | 5.00%   |
| Third Tone Sandhi Recognition (NCCU Corpus of Spoken Taiwan Mandarin)   | NAR↓   | 25.00%        | 34.38% | 65.62%     | 28.12%      | 62.50%          | 3.12%                | 28.12% | 56.25%   | 50.00%  |
| Third Tone Sandhi Recognition (NCCU Corpus of Spoken Taiwan Mandarin)   | IoU↑   | 0.0420        | 0.0530 | 0.0910     | 0.1460      | 0.0830          | 0.0650               | 0.1740 | 0.1430   | 0.0630  |
| L2 English Fluency (speechocean762-Ranking)                             | LLM-C↑ | 27.80%        | 22.80% | 50.00%     | 50.60%      | 51.90%          | 54.20%               | 30.00% | 50.60%   | 49.40%  |

Continued on next page

| Task   | Metric | Whisper-LLaMA | LTU-AS | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM | MU-LLaMA | GAMA-IT |
|--|--------|---------------|--------|------------|-------------|-----------------|----------------------|--------|----------|---------|
| L2 English Fluency (speechocean762-Scoring)      | NAR↓   | 13.59%        | 74.02% | 0.00%      | 19.52%      | 2.42%           | 0.40%                | 41.18% | 14.27%   | 2.02%   |
| L2 English Fluency (speechocean762-Scoring)      | PCC↑   | 0.0210        | 0.2580 | 0.0290     | 0.0370      | 0.0460          | -0.0860              | 0.0540 | 0.3110   | 0.2920  |
| L2 English Prosodic (speechocean762-Ranking)     | LLM-C↑ | 23.60%        | 36.70% | 46.70%     | 50.80%      | 36.10%          | 47.20%               | 45.30% | 52.50%   | 24.20%  |
| L2 English Prosodic (speechocean762-Scoring)     | NAR↓   | 27.86%        | 92.73% | 0.27%      | 10.77%      | 26.92%          | 2.69%                | 13.86% | 82.10%   | 17.77%  |
| L2 English Prosodic (speechocean762-Scoring)     | PCC↑   | 0.0900        | 0.3550 | 0.0430     | 0.0900      | 0.4930          | -0.0170              | 0.0940 | -0.0670  | 0.1330  |
| Prosody Naturalness (ProsAudit-Lexical)          | LLM-C↑ | 57.10%        | 40.90% | 51.40%     | 52.90%      | 32.00%          | 58.70%               | 68.70% | 74.90%   | 36.30%  |
| Prosody Naturalness (ProsAudit-Protosyntax)      | LLM-C↑ | 61.50%        | 41.60% | 53.80%     | 50.40%      | 36.30%          | 59.20%               | 59.20% | 75.20%   | 36.30%  |
| Spoof Detection (ASVspoof2015)                   | LLM-C↑ | 45.50%        | 23.00% | 56.00%     | 19.00%      | 13.00%          | 24.00%               | 23.00% | 73.00%   | 10.50%  |
| Spoof Detection (ASVspoof2017)                   | LLM-C↑ | 48.50%        | 11.00% | 68.50%     | 42.50%      | 28.00%          | 35.50%               | 27.50% | 66.00%   | 10.50%  |
| Deep Fake Voice Recognition (DEEP-VOICE)         | LLM-C↑ | 37.50%        | 26.00% | 51.80%     | 47.30%      | 49.50%          | 47.00%               | 47.50% | 42.00%   | 31.50%  |
| Enhancement Detection (LibriTTS-TestClean, WHAM) | LLM-C↑ | 52.00%        | 33.00% | 30.00%     | 50.50%      | 59.00%          | 58.00%               | 35.00% | 54.50%   | 45.50%  |
| Fraud Robocall Recognition (CallHome)            | LLM-C↑ | 73.30%        | 13.30% | 0.00%      | 100.00%     | 100.00%         | 100.00%              | 83.30% | 50.00%   | 100.00% |
| Fraud Robocall Recognition (Promo)               | LLM-C↑ | 47.40%        | 0.00%  | 0.00%      | 73.70%      | 42.10%          | 63.20%               | 47.40% | 47.40%   | 100.00% |
| Fraud Robocall Recognition (Robocall)            | LLM-C↑ | 94.90%        | 46.20% | 100.00%    | 20.50%      | 74.40%          | 61.50%               | 33.30% | 38.50%   | 0.00%   |
| HEAR Language Identification (VoxLingua107)      | LLM-C↑ | 88.68%        | 6.66%  | 10.48%     | 17.50%      | 26.02%          | 54.32%               | 4.70%  | 12.34%   | 3.10%   |
| Language Identification (VoxForge)               | LLM-C↑ | 96.00%        | 12.00% | 21.00%     | 8.50%       | 84.50%          | 93.00%               | 17.50% | 5.50%    | 0.00%   |
| Age Classification (Common Voice)                | LLM-C↑ | 12.80%        | 24.50% | 23.30%     | 20.80%      | 24.50%          | 36.50%               | 21.80% | 23.50%   | 15.50%  |
| Gender Recognition by Voice (Common Voice)       | LLM-C↑ | 4.50%         | 82.00% | 61.00%     | 57.00%      | 53.50%          | 97.50%               | 35.00% | 51.00%   | 64.50%  |

Continued on next page

| Task  | Metric | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM  | MU-LLaMA | GAMA-IT |
|---|--------|---------------|---------|------------|-------------|-----------------|----------------------|---------|----------|---------|
| HEAR Speaker Count Identification (LibriCount)    | NAR↓   | 32.80%        | 56.70%  | 45.30%     | 47.00%      | 10.70%          | 9.60%                | 37.90%  | 36.70%   | 24.80%  |
| HEAR Speaker Count Identification (LibriCount)    | ACC↑   | 20.39%        | 24.94%  | 16.64%     | 20.94%      | 12.21%          | 15.60%               | 11.11%  | 17.69%   | 17.42%  |
| Multi-Speaker Detection (VCTK)                    | LLM-C↑ | 49.50%        | 32.50%  | 36.00%     | 30.50%      | 50.00%          | 51.00%               | 30.50%  | 32.50%   | 42.00%  |
| Speaker Counting (LibriTTS-TestClean)             | LLM-C↑ | 18.00%        | 6.00%   | 14.50%     | 13.50%      | 17.50%          | 26.50%               | 14.00%  | 3.50%    | 14.50%  |
| Speaker Verification (LibriSpeech-TestClean)      | LLM-C↑ | 56.50%        | 28.50%  | 45.00%     | 58.00%      | 43.50%          | 63.00%               | 16.50%  | 36.00%   | 34.50%  |
| Speaker Verification (LibriSpeech-TestOther)      | LLM-C↑ | 45.00%        | 26.00%  | 51.00%     | 55.00%      | 49.00%          | 57.00%               | 20.00%  | 36.00%   | 33.00%  |
| Speaker Verification (VCTK)                       | LLM-C↑ | 55.00%        | 29.50%  | 54.00%     | 59.50%      | 44.00%          | 60.50%               | 23.00%  | 33.50%   | 36.50%  |
| SuperbSD (Libri2Mix-Test)                         | NAR↓   | 100.00%       | 100.00% | 100.00%    | 100.00%     | 100.00%         | 100.00%              | 100.00% | 100.00%  | 100.00% |
| SuperbSD (Libri2Mix-Test)                         | DER↓   | -             | -       | -          | -           | -               | -                    | -       | -        | -       |
| SuperbSV (SuperbHiddenSet)                        | LLM-C↑ | 53.00%        | 34.00%  | 45.00%     | 49.00%      | 49.00%          | 57.00%               | 48.00%  | 51.00%   | 39.00%  |
| Noise Detection (LJSpeech, MUSAN-Gaussian)        | LLM-C↑ | 45.50%        | 18.50%  | 50.00%     | 40.00%      | 42.00%          | 49.00%               | 45.00%  | 36.50%   | 31.50%  |
| Noise Detection (LJSpeech, MUSAN-Music)           | LLM-C↑ | 53.00%        | 23.00%  | 52.00%     | 47.00%      | 50.50%          | 47.50%               | 50.00%  | 31.50%   | 37.00%  |
| Noise Detection (LJSpeech, MUSAN-Noise)           | LLM-C↑ | 49.50%        | 18.00%  | 50.50%     | 57.00%      | 50.50%          | 59.50%               | 48.00%  | 48.00%   | 33.00%  |
| Noise Detection (LJSpeech, MUSAN-Speech)          | LLM-C↑ | 48.00%        | 22.00%  | 53.00%     | 36.00%      | 50.50%          | 47.00%               | 46.50%  | 50.00%   | 51.00%  |
| Noise Detection (VCTK, MUSAN-Music)               | LLM-C↑ | 47.50%        | 39.00%  | 45.50%     | 57.50%      | 42.00%          | 58.50%               | 56.00%  | 37.50%   | 39.00%  |
| Noise Detection (VCTK, MUSAN-Noise)               | LLM-C↑ | 45.50%        | 35.00%  | 45.50%     | 56.50%      | 44.50%          | 67.00%               | 53.50%  | 54.00%   | 33.50%  |
| Noise Detection (VCTK, MUSAN-Speech)              | LLM-C↑ | 46.50%        | 50.00%  | 43.50%     | 52.50%      | 45.50%          | 57.00%               | 46.50%  | 48.50%   | 56.50%  |
| Noise Detection (VCTK, MUSAN-Gaussian)            | LLM-C↑ | 53.00%        | 24.50%  | 49.00%     | 51.50%      | 45.50%          | 54.50%               | 52.50%  | 36.00%   | 21.00%  |
| Noise SNR Level Prediction (VCTK, MUSAN-Gaussian) | LLM-C↑ | 23.50%        | 15.00%  | 31.50%     | 21.00%      | 15.00%          | 22.00%               | 17.50%  | 21.50%   | 17.00%  |

Continued on next page

| Task  | Metric | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM  | MU-LLaMA | GAMA-IT |
|---|--------|---------------|---------|------------|-------------|-----------------|----------------------|---------|----------|---------|
| Noise SNR Level Prediction (VCTK, MUSAN-Music)            | LLM-C↑ | 21.50%        | 13.50%  | 24.50%     | 26.00%      | 14.50%          | 15.00%               | 17.00%  | 14.50%   | 10.00%  |
| Noise SNR Level Prediction (VCTK, MUSAN-Noise)            | LLM-C↑ | 30.00%        | 19.50%  | 10.50%     | 21.00%      | 10.50%          | 21.00%               | 16.50%  | 15.50%   | 15.00%  |
| Noise SNR Level Prediction (VCTK, MUSAN-Speech)           | LLM-C↑ | 24.00%        | 15.00%  | 16.00%     | 10.00%      | 17.50%          | 21.50%               | 18.00%  | 15.50%   | 13.50%  |
| Reverberation Detection (LJSpeech, RirsNoises-LargeRoom)  | LLM-C↑ | 46.50%        | 38.50%  | 45.50%     | 57.00%      | 54.50%          | 48.00%               | 47.00%  | 49.00%   | 27.50%  |
| Reverberation Detection (LJSpeech, RirsNoises-MediumRoom) | LLM-C↑ | 40.50%        | 25.50%  | 44.00%     | 51.00%      | 53.50%          | 50.50%               | 49.50%  | 48.00%   | 29.00%  |
| Reverberation Detection (LJSpeech, RirsNoises-SmallRoom)  | LLM-C↑ | 46.00%        | 36.50%  | 42.50%     | 47.50%      | 43.50%          | 51.00%               | 50.00%  | 44.50%   | 37.50%  |
| Reverberation Detection (VCTK, RirsNoises-LargeRoom)      | LLM-C↑ | 44.50%        | 37.50%  | 42.50%     | 48.00%      | 63.00%          | 46.00%               | 44.00%  | 45.00%   | 36.00%  |
| Reverberation Detection (VCTK, RirsNoises-MediumRoom)     | LLM-C↑ | 49.00%        | 32.00%  | 40.50%     | 45.50%      | 53.00%          | 51.50%               | 44.00%  | 47.00%   | 27.00%  |
| Reverberation Detection (VCTK, RirsNoises-SmallRoom)      | LLM-C↑ | 48.00%        | 39.00%  | 42.50%     | 44.00%      | 47.50%          | 51.50%               | 45.50%  | 43.50%   | 47.00%  |
| N-Best Correction (Librispeech-TestOther)                 | LLM-C↑ | 45.00%        | 30.00%  | 31.60%     | 36.20%      | 42.60%          | 41.80%               | 39.20%  | 39.60%   | 35.40%  |
| AAVE Speech Recognition (CORAAL)                          | WER↓   | 21.70%        | 97.30%  | 24.00%     | 31.60%      | 96.80%          | 38.60%               | 34.90%  | 136.50%  | 102.90% |
| Code-switch Speech Recognition (NTUML2021)                | MER↓   | 424.58%       | 215.51% | 293.50%    | 185.27%     | 165.54%         | 116.88%              | 172.03% | 130.18%  | 193.47% |
| Code-Switching Speech Recognition (ASCEND)                | NAR↓   | 68.20%        | 12.20%  | 5.40%      | 88.80%      | 1.60%           | 0.00%                | 30.40%  | 69.40%   | 14.60%  |
| Code-Switching Speech Recognition (ASCEND)                | ACC↑   | 28.30%        | 50.57%  | 58.14%     | 14.29%      | 44.31%          | 60.80%               | 10.92%  | 13.73%   | 9.84%   |
| Multi-Lingual Speech Recognition (MLS-de)                 | WER↓   | 73.80%        | 132.35% | 34.46%     | 25.21%      | 79.60%          | 24.96%               | 49.56%  | 99.37%   | 105.81% |
| Multi-Lingual Speech Recognition (MLS-en)                 | WER↓   | 65.03%        | 100.29% | 9.37%      | 9.47%       | 27.97%          | 17.47%               | 17.23%  | 96.03%   | 97.09%  |

Continued on next page



| Task  | Metric | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM  | MU-LLaMA | GAMA-IT  |
|---|--------|---------------|---------|------------|-------------|-----------------|----------------------|---------|----------|----------|
| Multi-Lingual Speech Recognition (MLS-es)       | WER↓   | 72.49%        | 123.26% | 23.45%     | 16.10%      | 58.82%          | 18.19%               | 44.90%  | 103.19%  | 101.75%  |
| Multi-Lingual Speech Recognition (MLS-fr)       | WER↓   | 71.09%        | 133.86% | 26.52%     | 21.27%      | 46.27%          | 19.01%               | 42.35%  | 102.25%  | 102.01%  |
| Multi-Lingual Speech Recognition (MLS-it)       | WER↓   | 85.38%        | 125.39% | 39.11%     | 31.54%      | 56.55%          | 33.16%               | 46.14%  | 102.72%  | 108.84%  |
| Multi-Lingual Speech Recognition (MLS-nl)       | WER↓   | 74.96%        | 120.29% | 32.11%     | 29.37%      | 140.01%         | 42.94%               | 55.56%  | 102.05%  | 101.22%  |
| Multi-Lingual Speech Recognition (MLS-pl)       | WER↓   | 76.79%        | 139.04% | 73.83%     | 51.53%      | 278.19%         | 101.41%              | 90.11%  | 100.23%  | 105.43%  |
| Multi-Lingual Speech Recognition (MLS-pt)       | WER↓   | 73.63%        | 136.66% | 31.88%     | 25.07%      | 117.33%         | 24.72%               | 46.26%  | 99.71%   | 103.48%  |
| PTBR Speech Recognition (Common Voice)          | WER↓   | 45.75%        | 248.05% | 69.88%     | 73.91%      | 233.85%         | 38.57%               | 62.93%  | 639.12%  | 270.64%  |
| SuperbASR (LibriSpeech-TestClean)               | WER↓   | 33.96%        | 96.28%  | 15.11%     | 2.79%       | 69.35%          | 36.70%               | 6.87%   | 103.67%  | 116.66%  |
| SuperbASR (LibriSpeech-TestOther)               | WER↓   | 42.14%        | 91.89%  | 11.44%     | 4.31%       | 79.53%          | 40.28%               | 9.17%   | 111.03%  | 130.00%  |
| SuperbOOD Asr Ar (Common Voice)                 | WER↓   | 51.04%        | 245.85% | 216.51%    | 178.02%     | 289.25%         | 178.21%              | 149.15% | 225.28%  | 504.53%  |
| SuperbOOD Asr Es (Common Voice)                 | WER↓   | 10.93%        | 150.00% | 98.96%     | 99.22%      | 100.83%         | 75.08%               | 99.38%  | 141.68%  | 303.28%  |
| SuperbOOD Asr Spon (CHIME6-Test)                | WER↓   | 65.71%        | 122.86% | 61.11%     | 62.80%      | 92.17%          | 80.21%               | 136.24% | 137.25%  | 266.72%  |
| SuperbOOD Asr Zh (Common Voice)                 | CER↓   | 29.67%        | 508.02% | 609.28%    | 310.33%     | 449.32%         | 270.63%              | 445.36% | 435.33%  | 1097.41% |
| Target Speaker ASR (AMI-test)                   | WER↓   | 143.02%       | 133.86% | 273.90%    | 187.15%     | 207.01%         | 132.13%              | 266.26% | 108.32%  | 140.37%  |
| Multi-Speaker Detection (LibriSpeech-TestClean) | LLM-C↑ | 49.00%        | 31.50%  | 16.50%     | 21.50%      | 50.50%          | 52.00%               | 20.50%  | 33.50%   | 48.50%   |
| Speech Command Recognition (AudioMNIST)         | NAR↓   | 4.93%         | 76.80%  | 2.93%      | 0.00%       | 0.93%           | 0.00%                | 0.53%   | 5.20%    | 75.87%   |
| Speech Command Recognition (AudioMNIST)         | ACC↑   | 0.8890        | 0.8330  | 0.9670     | 0.7410      | 0.9620          | 0.7710               | 0.9340  | 0.1000   | 0.2430   |
| Speech Text Matching (LibriSpeech-TestClean)    | LLM-C↑ | 87.00%        | 32.00%  | 58.00%     | 57.00%      | 68.00%          | 84.50%               | 52.00%  | 50.00%   | 18.50%   |

Continued on next page

| Task  | Metric | Whisper-LLaMA | LTU-AS  | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM  | MU-LLaMA | GAMA-IT |
|---|--------|---------------|---------|------------|-------------|-----------------|----------------------|---------|----------|---------|
| Speech Text Matching (LibriSpeech-TestOther)  | LLM-C↑ | 82.00%        | 34.00%  | 57.00%     | 56.50%      | 72.00%          | 87.50%               | 51.50%  | 47.50%   | 18.50%  |
| Speech Text Matching (LJSpeech)               | LLM-C↑ | 82.50%        | 32.00%  | 59.50%     | 55.50%      | 64.50%          | 76.00%               | 50.50%  | 39.50%   | 16.00%  |
| Spoken Term Detection (LibriSpeech-TestClean) | LLM-C↑ | 77.00%        | 29.50%  | 54.50%     | 53.50%      | 52.50%          | 62.00%               | 48.50%  | 32.00%   | 39.50%  |
| Spoken Term Detection (LibriSpeech-TestOther) | LLM-C↑ | 77.00%        | 27.50%  | 50.00%     | 50.00%      | 48.50%          | 66.00%               | 47.00%  | 30.50%   | 39.50%  |
| Spoken Term Detection (LJSpeech)              | LLM-C↑ | 83.50%        | 28.50%  | 53.00%     | 51.50%      | 47.00%          | 79.50%               | 50.00%  | 28.50%   | 31.00%  |
| SuperbKS (Speech Commands V1-Test)            | LLM-C↑ | 75.50%        | 4.50%   | 30.50%     | 31.50%      | 30.00%          | 59.00%               | 47.50%  | 27.50%   | 34.00%  |
| SuperbQbE (Quesst14-Eval)                     | LLM-C↑ | 48.50%        | 33.50%  | 49.00%     | 51.50%      | 37.50%          | 53.50%               | 49.00%  | 52.00%   | 31.00%  |
| Stuttering Detection (SEP28k)                 | LLM-C↑ | 49.50%        | 45.40%  | 50.50%     | 50.30%      | 50.20%          | 52.60%               | 58.90%  | 49.80%   | 48.50%  |
| Voice Disorder Classification (VOICED)        | LLM-C↑ | 16.30%        | 13.50%  | 13.50%     | 17.30%      | 19.70%          | 17.30%               | 20.20%  | 22.10%   | 13.50%  |
| Conversation Matching                         | LLM-C↑ | 77.80%        | 9.30%   | 57.40%     | 37.00%      | 51.90%          | 61.10%               | 61.10%  | 1.90%    | 18.50%  |
| Dialogue Act Classification (SLUE-HVB)        | LLM-C↑ | 21.80%        | 18.90%  | 22.80%     | 3.30%       | 9.40%           | 20.60%               | 17.30%  | 4.90%    | 4.10%   |
| Dialogue Act Classification (Daily Talk)      | LLM-C↑ | 29.00%        | 14.50%  | 37.00%     | 44.00%      | 43.50%          | 30.50%               | 42.00%  | 20.50%   | 11.00%  |
| Dialogue Act Pairing (Daily Talk)             | LLM-C↑ | 51.00%        | 16.00%  | 53.50%     | 48.00%      | 43.00%          | 48.50%               | 39.50%  | 50.00%   | 39.00%  |
| SuperbIC (SLURP)                              | LLM-C↑ | 33.00%        | 7.50%   | 13.50%     | 3.00%       | 31.00%          | 26.50%               | 20.50%  | 4.00%    | 2.00%   |
| Intent Classification (SLURP-Action)          | LLM-C↑ | 24.50%        | 12.50%  | 16.50%     | 16.00%      | 27.50%          | 67.50%               | 33.00%  | 7.00%    | 10.00%  |
| Intent Classification (SLURP-Intent)          | LLM-C↑ | 36.50%        | 12.00%  | 21.00%     | 13.50%      | 48.50%          | 43.00%               | 41.50%  | 6.00%    | 6.50%   |
| Named Entity Localization (SLUE-VoxPopuli)    | NAR↓   | 100.00%       | 100.00% | 100.00%    | 100.00%     | 100.00%         | 100.00%              | 100.00% | 100.00%  | 100.00% |
| Named Entity Localization (SLUE-VoxPopuli)    | F1↑    | -             | -       | -          | -           | -               | -                    | -       | -        | -       |
| Named Entity Recognition (SLUE-VoxPopuli)     | NAR↓   | 43.48%        | 84.78%  | 67.39%     | 55.43%      | 44.57%          | 64.13%               | 68.48%  | 89.13%   | 90.22%  |

*Continued on next page*

| Task  | Metric          | Whisper-LLaMA | LTU-AS | SALMONN 7B | SALMONN 13B | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM  | MU-LLaMA | GAMA-IT |
|---|-----------------|---------------|--------|------------|-------------|-----------------|----------------------|---------|----------|---------|
| Named Entity Recognition (SLUE-VoxPopuli)                             | IoU↑            | 0.0190        | 0.2640 | 0.2120     | 0.1180      | 0.0290          | 0.1190               | 0.6930  | 0.0500   | 0.0000  |
| Sarcasm Detection (Mustard)   | LLM-C↑          | 38.50%        | 36.50% | 41.50%     | 46.00%      | 44.00%          | 48.00%               | 47.50%  | 45.50%   | 36.00%  |
| Semantic Textual Similarity (SpokenSTS)                               | LLM-C↑          | 50.40%        | 40.80% | 49.60%     | 50.40%      | 22.00%          | 49.20%               | 43.60%  | 54.00%   | 21.20%  |
| Speech Sentiment Analysis (MELD)                                      | LLM-C↑          | 39.00%        | 31.00% | 33.50%     | 36.00%      | 38.50%          | 49.00%               | 42.50%  | 28.50%   | 29.50%  |
| Spoken Digit Arithmetic (AudioMNIST)                                  | LLM-C↑          | 43.50%        | 13.00% | 31.00%     | 33.00%      | 51.50%          | 43.50%               | 41.50%  | 26.50%   | 16.50%  |
| SuperbSF (AudioSnips-Test)  | NAR↓            | 40.50%        | 82.50% | 44.00%     | 70.00%      | 41.50%          | 58.00%               | 63.50%  | 95.00%   | 92.00%  |
| SuperbSF (AudioSnips-Test)  | Slot Type F1↑   | 0.8570        | 0.9020 | 0.9230     | 0.9210      | 0.8650          | 0.9270               | 0.8960  | 0.7770   | 0.9020  |
| SuperbSF (AudioSnips-Test)  | Slot Value CER↓ | 0.3540        | 0.8500 | 0.2370     | 0.2970      | 0.4600          | 0.2390               | 0.5510  | 0.6690   | 0.4160  |
| Code-Switching Semantic Grammar Acceptability Comparison (CSZS-zh-en) | LLM-C↑          | 51.50%        | 37.50% | 50.00%     | 49.50%      | 45.50%          | 49.50%               | 23.00%  | 45.50%   | 32.00%  |
| Nonce Word Detection (sWUGGY)   | LLM-C↑          | 63.00%        | 25.60% | 55.00%     | 58.10%      | 32.20%          | 61.80%               | 63.50%  | 42.50%   | 19.90%  |
| PoS Estimation (LibriTTS)   | POS↓            | 2.5722        | 1.2792 | 2.9675     | 1.1990      | 3.0620          | 2.3091               | 3.1126  | 1.5327   | 1.9635  |
| PoS Estimation with transcription (LibriTTS)                          | POS↓            | 0.9081        | 1.3507 | 1.7583     | 1.0601      | 2.1580          | 1.7514               | 2.0463  | 1.2394   | 1.3911  |
| Sentence Grammar Acceptability (sBLIMP)                               | LLM-C↑          | 68.50%        | 7.30%  | 44.00%     | 62.30%      | 50.80%          | 73.80%               | 61.30%  | 72.80%   | 10.50%  |
| SuperbST (CoVoST2-Test)   | Sacre Bleu↑     | 17.3720       | 0.1070 | 18.7884    | 16.7835     | 6.9846          | 23.3458              | 21.7423 | 0.0566   | 0.0219  |

## D HEAR & MARBLE EVALUATION RESULTS

Tables 7 and 8 show the evaluation results for the HEAR and MARBLE tasks, respectively, as collected in Dynamic-SUPERB *Phase-2*. Importantly, we replaced the datasets for some tasks to address licensing issues and reduced their size to enable more efficient inference. Furthermore, we removed certain tasks because they overlap with existing SUPERB tasks. Consequently, these results are not directly comparable to those evaluated on the original HEAR and MARBLE.

Table 7: HEAR (audio) Results. The SpeechCommands-5hr task results are omitted as they duplicate the SUPERB Keyword Spotting (KS) task, presented in Table 2. *LLM-C* stands for LLM Classification, *NAR* stands for Not Applicable Rate, *TER* stands for Token Error Rate, and *ACC* stands for Accuracy.

| Task  | Metric                               | Whisper-LLaMA    | LTU-AS           | SALMONN 7B        | SALMONN 13B      | Qwen-Audio-Chat  | Qwen2-Audio-7B-Inst. | WavLLM           | MU-LLaMA         | GAMA-IT          |
|---|--------------------------------------|------------------|------------------|-------------------|------------------|------------------|----------------------|------------------|------------------|------------------|
| Beehive States Classification - Beehive States                              | LLM-C $\uparrow$                     | 43.00%           | 35.00%           | 54.00%            | 42.50%           | 19.00%           | 42.50%               | 37.00%           | 37.50%           | 19.00%           |
| Emotion Recognition - CREMAD  | LLM-C $\uparrow$                     | 6.70%            | 9.30%            | 18.80%            | 18.10%           | 62.50%           | 61.10%               | 25.60%           | 12.40%           | 13.20%           |
| Environmental Sound Classification - ESC-50                                 | LLM-C $\uparrow$                     | 4.80%            | 25.30%           | 61.00%            | 73.40%           | 68.80%           | 43.90%               | 9.80%            | 0.40%            | 35.30%           |
| Language Identification - VoxLingua107 Top10                                | LLM-C $\uparrow$                     | 92.18%           | 1.95%            | 6.28%             | 11.01%           | 18.00%           | 36.11%               | 1.85%            | 0.00%            | 0.10%            |
| Music Genre Classification - ISMIR04  | LLM-C $\uparrow$                     | 15.58%           | 7.04%            | 31.16%            | 44.22%           | 33.67%           | 37.69%               | 16.58%           | 27.64%           | 2.51%            |
| Music Speech Classification - MAESTRO-LibriSpeech                           | LLM-C $\uparrow$                     | 72.50%           | 77.50%           | 91.50%            | 99.50%           | 98.00%           | 89.50%               | 59.50%           | 51.50%           | 88.00%           |
| Music Transcription - MAESTRO   | NAR $\downarrow$<br>TER $\downarrow$ | 92.43%<br>1.0065 | 100.00%<br>-     | 99.46%<br>36.0000 | 100.00%<br>-     | 96.76%<br>8.3438 | 96.76%<br>1.6875     | 96.76%<br>6.2500 | 100.00%<br>-     | 100.00%<br>-     |
| Percussion Instruments Classification - Beijing Opera Percussion Instrument | LLM-C $\uparrow$                     | 3.39%            | 4.24%            | 10.17%            | 20.76%           | 16.10%           | 27.12%               | 24.15%           | 22.03%           | 11.44%           |
| Percussion Instruments Stroke Classification - Mridangam Stroke             | LLM-C $\uparrow$                     | 0.10%            | 3.80%            | 7.60%             | 4.90%            | 11.50%           | 10.80%               | 13.70%           | 10.30%           | 0.10%            |
| Percussion Instruments Tonic Classification - Mridangam Stroke              | LLM-C $\uparrow$                     | 2.00%            | 2.50%            | 12.70%            | 13.40%           | 11.30%           | 11.80%               | 12.80%           | 12.50%           | 5.70%            |
| Sound Event Detection - DCASE2016 Task2                                     | LLM-C $\uparrow$                     | 0.00%            | 0.00%            | 21.43%            | 35.71%           | 7.14%            | 14.29%               | 21.43%           | 21.43%           | 7.14%            |
| Speaker Count Identification - LibriCount                                   | NAR $\downarrow$<br>ACC $\uparrow$   | 30.40%<br>17.53% | 42.10%<br>16.93% | 35.90%<br>13.26%  | 42.50%<br>17.22% | 3.50%<br>10.26%  | 0.40%<br>13.76%      | 22.00%<br>5.51%  | 23.20%<br>13.80% | 26.20%<br>17.48% |
| Vocal Imitation Classification - Vocal Imitations                           | LLM-C $\uparrow$                     | 15.17%           | 5.50%            | 1.17%             | 1.83%            | 18.67%           | 11.17%               | 17.50%           | 3.83%            | 5.83%            |

Table 8: MARBLE (music) Results. *LLM-C* stands for LLM Classification, *NAR* stands for Not Applicable Rate, and *MSE* stands for Mean Squared Error.

| Task                                    | Metric                             | Whisper-LLaMA | LTU-AS       | SALMONN 7B   | SALMONN 13B  | Qwen-Audio-Chat | Qwen2-Audio-7B-Inst. | WavLLM       | MU-LLaMA     | GAMA-IT      |
|---|------------------------------------|---------------|--------------|--------------|--------------|-----------------|----------------------|--------------|--------------|--------------|
| Key Detection - Giantsteps (Key)        | LLM-C $\uparrow$                   | 3.00%         | 0.50%        | 1.50%        | 1.00%        | 3.00%           | 17.00%               | 2.00%        | 0.50%        | 0.00%        |
| Music Tagging - MTG-Jamendo             | LLM-C $\uparrow$                   | 0.10%         | 0.30%        | 0.00%        | 0.00%        | 2.90%           | 0.30%                | 0.10%        | 0.50%        | 0.10%        |
| Music Tagging - MagnaTagATune           | LLM-C $\uparrow$                   | 5.50%         | 1.00%        | 3.50%        | 9.00%        | 4.50%           | 1.50%                | 1.00%        | 3.50%        | 7.50%        |
| Genre Classification - MTG-Jamendo      | LLM-C $\uparrow$                   | 0.00%         | 1.60%        | 0.10%        | 0.20%        | 3.20%           | 0.10%                | 0.20%        | 0.20%        | 0.00%        |
| Emotion Detection - MTG-Jamendo         | LLM-C $\uparrow$                   | 0.00%         | 1.50%        | 0.00%        | 0.00%        | 2.70%           | 0.80%                | 0.80%        | 0.80%        | 0.10%        |
| Beat Tracking - ASAP                    | NAR $\downarrow$<br>MSE $\uparrow$ | 100.00%<br>-  | 100.00%<br>- | 100.00%<br>- | 100.00%<br>- | 100.00%<br>-    | 100.00%<br>-         | 100.00%<br>- | 100.00%<br>- | 100.00%<br>- |
| Vocal Technique Detection - VocalSet    | LLM-C $\uparrow$                   | 8.00%         | 1.00%        | 11.00%       | 6.50%        | 6.50%           | 15.50%               | 8.00%        | 8.50%        | 2.50%        |
| Instrument Classification - MTG-Jamendo | LLM-C $\uparrow$                   | 0.40%         | 2.30%        | 0.00%        | 0.00%        | 6.40%           | 1.40%                | 0.20%        | 0.60%        | 0.20%        |

## E LLM-BASED EVALUATION

Here, we list the prompt templates used in the evaluation. Figure 6 illustrates the prompt for classification task evaluation, while Figure 7 shows the prompt for regression task evaluation. In addition to GPT-4o, we also employed open-source LLMs to ensure better reproducibility. Tables 9 and 10 present the performance evaluations of the proposed prompts on classification and regression tasks, respectively, comparing their correlation with a human-labeled evaluation set. We found that GPT-4o and LLaMA-3.1-70B-Instruct achieved very high accuracy in both classification and regression tasks. This demonstrates that our evaluation using GPT-4o is reliable and that researchers can leverage LLaMA-3.1-70B-Instruct to conduct fully reproducible evaluations.

Table 9: Performance evaluation of the proposed post-processing prompt on **classification** tasks using a human-labeled evaluation set. Accuracy is determined by agreement between the model output and the corresponding human label.

| Model Name             | Accuracy      | Correct Count  |
|------------------------|---------------|----------------|
| gpt-4o                 | <b>98.67%</b> | <b>148/150</b> |
| gpt-4o-mini            | 98.00%        | 147/150        |
| gpt-o1-preview         | 98.67%        | 148/150        |
| gpt-o1-mini            | <b>99.33%</b> | <b>149/150</b> |
| llama-3.1-70B-Instruct | <b>96.67%</b> | <b>145/150</b> |
| llama-3.1-8B-Instruct  | 92.67%        | 139/150        |
| llama-3.2-3B-Instruct  | 86.67%        | 130/150        |
| llama-3.2-1B-Instruct  | 61.33%        | 92/150         |

Table 10: Performance evaluation of the proposed post-processing prompt on **regression** tasks using a human-labeled evaluation set. Accuracy is determined by exact matches between the post-processed model output and the corresponding human-generated post-processed output.

| Model Name             | Accuracy      | Correct Count  |
|------------------------|---------------|----------------|
| gpt-4o                 | <b>92.00%</b> | <b>138/150</b> |
| gpt-4o-mini            | 86.67%        | 130/150        |
| gpt-o1-preview         | 86.00%        | 129/150        |
| gpt-o1-mini            | 89.33%        | 134/150        |
| llama-3.1-70B-Instruct | <b>90.00%</b> | <b>135/150</b> |
| llama-3.1-8B-Instruct  | 82.00%        | 123/150        |
| llama-3.2-3B-Instruct  | 61.33%        | 92/150         |
| llama-3.2-1B-Instruct  | 39.33%        | 59/150         |

**Classification Post-processing Prompt:**

You will be given a question, a corresponding correct answer(s), and a response from a model. The model’s response is a reply to the question. Your task is to judge if the “Model’s Response” aligns with the “Ground Truth Answer” based on the “Question.”

Please strictly follow the guidelines below:

- Briefly explain the reasons for your judgment.
- Answer with the format “Result: <YES or NO>” at the end.
- Output “YES” if the response aligns with the ground truth answer; output “NO” if the response does not match the ground truth answer, selects incorrect or irrelevant options, or provides more answers than required.
- The questions would be single-choice or multi-choice:  
For single-choice questions, the model’s response should contain one and only one answer. If the model’s response selects more than one answer or does not clearly indicate a single answer, you should mark it as incorrect and output “NO.” For multi-choice questions, the model’s response must exactly match all applicable correct choices. If the model’s response selects too many, too few, or any incorrect answers, you should mark it as incorrect and output “NO.”
- Since the question is short answer, the model’s response does not need to mention the content of the question. You only need to check if the model’s response has the same meaning as the ground truth answer(s).

Input Format:

Question: {instruction}

Ground Truth Answer: {label}

Model’s Response: {response}

Figure 6: The prompt used for post-processing the **classification** tasks inference outputs in the study.

**Regression Post-processing Prompt:**

You will be provided with an "instruction," a "ground-truth label," and a "model output." Your job is to extract the value from the given "model output" and return the post-processed "model output". Provide your output with just the post-processed output without an explanation.

**[Task Description]**

- Analyze the "instruction" and "ground-truth label" to understand the desired output format.
- Post-process the "model output" so that it matches the format of the "ground-truth label" and is numerically comparable to it.

**[Guidelines]**

- If the "model output" can be adjusted to match the format of the "ground-truth label," return only the post-processed result.
- If the "model output" is already in an appropriate format, return it as is without changing it.
- If it can not be adjusted into the correct format, does not clearly indicate any of the candidates, is ambiguous, or implies more than one answer, return "N/A."
- Do not confuse the "ground-truth label" with "model output"; you are only supposed to use the "ground-truth label" to understand the desired output format.
- Do not tamper with the original content of "model output"; you are only supposed to post-process the format of "model output" so it matches the format of "ground-truth label" and can be compared or evaluated against each other directly.
- The post-processed output should only be a number in numeric form.

Below you have three examples. Please consider their patterns to understand your task better:

**[Example 1 - Inputs]**

- "instruction": "Please identify the total times of code-switching in this wavefile."
- "ground-truth label": 8
- "model output": "There are no instances of code-switching in this wavefile; it is purely Mandarin speech."
- Your response: 0

**[Example 2 - Inputs]**

- "instruction": ""Listen to the audio and assess the clarity and accuracy of pronunciation by considering precise phonology, pronunciation errors, and overall comprehensibility. Use an Arabic numeral from 0 to 10 for your answer.""
- "ground-truth label": 3
- "model output": "The audio is of a man speaking, in a Chinese, saying, "I has to find a different back door". The accuracy of the pronunciation is 90%."
- Your response: 9

**[Example 3 - Inputs]**

- "instruction": "Please predict the spatial distance(0.5m-4.5m) based on the given audio."
- "ground-truth label": 2.1196093634329904
- "model output": "Based on the given audio, the spatial distance is likely between 0.5m and 4.5m."
- Your response: N/A

**[Inputs]**

- "instruction": {instruction}
- "ground-truth label": {label}
- "model output": {model\_output}
- Your response:

Figure 7: The prompt used for post-processing the **regression** tasks inference outputs in the study.

## F CALL FOR TASKS

Here we provide material from the call for tasks. Figure 8 shows a task proposal that we ask each contributor to initiate. Task proposers need to provide a high-level description of the task, explain its importance and challenges, and list the datasets they plan to use along with the dataset licenses. Figure 9 presents the standard README format for each task. In the README, task proposers include the task introduction, its challenges, the dataset, and the evaluation metric. Importantly, we

also ask them to provide a table of state-of-the-art performance on the task, even if not achieved by universal spoken language models. These results can serve as a reference for researchers to better understand the task's difficulty and how it has developed so far. Then, Figure 10 shows the JSON format we require to record task information. For each task, our evaluation pipeline uses this JSON file as input to automatically download data from Huggingface, formatting it for subsequent use.

## Part-of-Speech Estimation

This task is to tag Part-of-Speech (POS) tags in an utterance. This task is significant because it tests the model's proficiency in understanding and processing the syntactic structure of spoken language. Accurate POS tagging from acoustic data can enhance various applications such as speech recognition (ASR), spoken language understanding (SLU), and machine translation (MT).

The task is part of the metatask [#140](#).

This task can provide greater insights when combined with our other task, [#142](#), which focuses on word-boundary detection. If the model performs poorly in both [#142](#) and this task, we can conclude that the model's inability to perform POS tagging is due to its failure to learn word boundaries. If the model performs well in task [#142](#) but poorly in this task, we can infer that the model lacks syntactic ability.

### Task Objective

- Accurately tag a word with POS tags from audio input.

### Datasets and Evaluation Metric

- Dataset
  - We plan to use LibriTTS test-other split [\[dataset\]](#) [\[paper\]](#)
  - We chose this dataset because the speech is split in sentence breaks. The dataset includes audio files and transcriptions.
  - For simplicity, we select sentences with 10~15 words. The dataset is under the CC BY 4.0 license.
- Preprocessing
  - Use **Stanza** to attain word-POS tag pair from transcriptions.
    - We use Stanza which uses Universal POS tag set, [\[universal POS\]](#).
    - Tag list: adjective (ADJ), adposition (ADP), adverb (ADV), auxiliary (AUX), coordinating conjunction (CCONJ), determiner (DET), interjection (INTJ), noun (NOUN), numeral (NUM), particle (PART), pronoun (PRON), proper noun (PROPN), subordinating conjunction (SCONJ), verb (VERB)
    - The reported performance on POS tagging is around 96.70% to 98.28% on official test sets for POS tagging [\[stanza\\_performance\]](#).

-Evaluation

- WER: word error rate of transcriptions
- POS Accuracy: accuracy of POS tags in an utterance
- Per POS Tag Accuracy: accuracy for each individual POS tag
- Sample instructions
  - For each word in the given audio utterance, determine and assign the corresponding Part-of-Speech (POS) tag. The list of tags includes ADJ (adjective), ADP (adposition), ADV (adverb), AUX (auxiliary), CCONJ (coordinating conjunction), DET (determiner), INTJ (interjection), NOUN (noun), NUM (numeral), PART (particle), PRON (pronoun), PROPN (proper noun), SCONJ (subordinating conjunction), and VERB (verb). Please write the transcribed utterance first, and then POS tags followed by a separate '/'. For example, when you listen to an audio of 'The quick brown fox', please generate the tags in the form of (transcription): The quick brown fox / (POS): DET ADJ ADJ NOUN.
  - For each word in the given audio clip, assign the appropriate Part-of-Speech (POS) tag. The available tags include ADJ (adjective), ADP (adposition), ADV (adverb), AUX (auxiliary), CCONJ (coordinating conjunction), DET (determiner), INTJ (interjection), NOUN (noun), NUM (numeral), PART (particle), PRON (pronoun), PROPN (proper noun), SCONJ (subordinating conjunction), and VERB (verb). Begin with the transcribed sentence, followed by the POS tags separated by a '/'. For instance, if the transcription is "She didn't eat breakfast," it should be formatted as: Transcription: She didn't eat breakfast / POS: PRON AUX VERB NOUN.

Figure 8: An example of a task proposal from a task contributor.

## Part of Speech Estimation

The Part of Speech (PoS) Estimation task aims to accurately tag words with their respective PoS tags from audio input. This task is crucial as it evaluates the model's capability in understanding and processing the syntactic structure of spoken language. Accurate PoS tagging from acoustic data can significantly enhance various applications such as Automatic Speech Recognition (ASR), Spoken Language Understanding (SLU), and Machine Translation (MT).

### Task Objective

The primary objective is to assess whether benchmark models can accurately identify the PoS of words in speech utterances. The task is divided into two subtasks:

- PoS estimation: The model is asked to generate PoS tags from the audio.
- PoS estimation with transcription: The model is asked to generate both transcriptions and PoS tags from the audio.

### Dataset

The dataset is derived from the LibriTTS corpus, a multi-speaker English corpus [\[LibriTTS-openSLR\]](#) [\[LibriTTS-HuggingFace\]](#). To suit the purpose of the benchmark task, we went through several preprocessing steps as follows: (1) Filter out sentences with quotes (ex. 'Yes', he said.) (2) Remove punctuations (3) Select sentences that include 3 to 15 words. We uniformly selected the sentences for each number of words. Consequently, the dataset includes 890 samples. Further details regarding the preprocessing steps, please refer to this repository [\[redacted\]](#)

The PoS Estimation task includes two variants: each subtask, PoS estimation and PoS estimation with transcription, corresponds to the PoSEstimation\_LibriTTS\_PoS dataset and the PoSEstimation\_LibriTTS\_sentPoS dataset, respectively. In PoSEstimation\_LibriTTS\_PoS, the model is asked to generate PoS tags only from the audio, and for PoSEstimation\_LibriTTS\_sentPoS, the model is asked to generate both transcriptions and PoS tags from the audio. Therefore, while the audio samples are shared between the two datasets, the 'instruction' and 'label' columns differ.

For more detailed information, please refer to the dataset directory : [\[redacted\]](#)

- version: 26bee87a4308883aaa21373f705d36b9d783d95e : [\[redacted\]](#)
- version: e31ef0a0b366c66d859542536e9ac5a8de92dac4

### Evaluation Metric

- Word Error Rate (WER): PoS tags are evaluated by using WER metric.
- WER with transcriptions: For the PoS estimation with transcription subtask, certain preprocessing is necessary to exclude the transcriptions, which are expected to be separated by '/'.

### State-of-the-art performance

To the best of our knowledge, LibriTTS has not been utilized for Part-of-Speech tagging tasks. However, PoS tagging tasks in English have been reported to show performance between 96.20% and 99.15% with conventional PoS taggers (<https://stanfordnlp.github.io/stanza/performance.html>).

### Reference

[1] Zen, H., Dang, V., Clark, R.A., Zhang, Y., Weiss, R.J., Jia, Y., Chen, Z., & Wu, Y. (2019). LibriTTS: A Corpus Derived from LibriSpeech for Text-to-Speech. Interspeech.

Figure 9: An example of the README file required in the task contribution.

```

1  {
2      "name": "PoSEstimation_LibriTTS_sentPoS",
3      "description": "Generate correct part-of-speech tags. Also evaluate the performance on transcriptions as a side task. We use word error rate metric for both tasks.",
4      "keywords": "linguistics",
5      "metrics": [
6          "word error rate (PoS)",
7          "word error rate (transcription)"
8      ],
9      "path": "DynamicSuperb/PoS_Estimation_LibriTTS_PoS_with_transcription",
10     "version": "0d2b8c9ad06d9558b257d4936c6d2b9641016301"
11 }

```

Figure 10: An example of the JSON file required in the task contribution.



## G AUDIO DURATION DISTRIBUTION

Figure 11 shows the distribution of all audio files in Dynamic-SUPERB *Phase-2*. Only about 6.5% of the audio files have a duration longer than 30 seconds, which reflects an inherent limitation of models that incorporate Whisper in their framework.

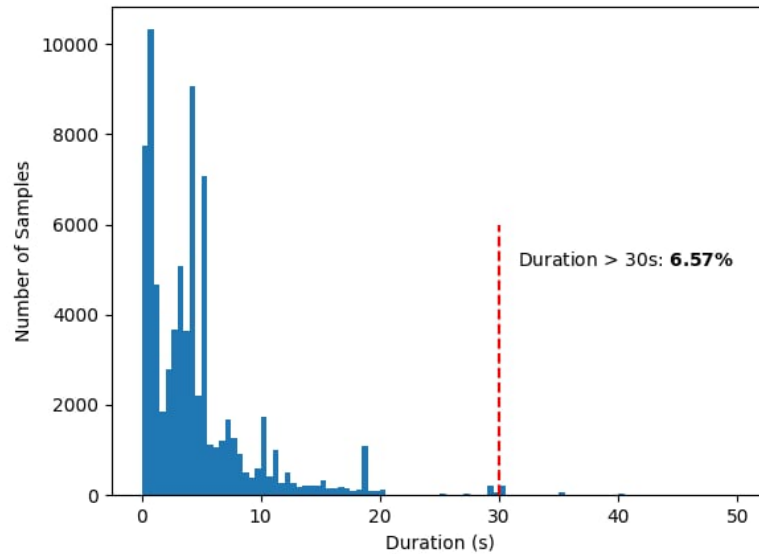


Figure 11: Duration distribution of all audios in Dynamic-SUPERB *Phase-2*.