

SEA-SpeechBench: A Large-Scale Multitask Benchmark for Speech Understanding Across Southeast Asia

Anonymous ACL submission

Abstract

The rapid advancement of audio and multi-modal large language models has unlocked transformative speech understanding capabilities, yet evaluation frameworks remain predominantly English-centric, leaving Southeast Asian (SEA) languages critically underrepresented. We introduce SEA-SPEECHBENCH, the first large-scale multitask benchmark that evaluates speech understanding in 11 SEA languages through more than 97,000 samples and 597 hours of curated audio data. Our benchmark comprises 9 diverse tasks across 3 categories: speech processing (automatic speech recognition, speech translation, spoken question answering), paralinguistic analysis (emotion, gender, age, speaker recognition), and temporal understanding, a novel dimension featuring timestamped content queries and temporal localization within extended audio sequences up to 3 minutes. We implement multilingual prompting in both native SEA languages and English to reflect user interactions with audio-language models. Evaluation of leading open-source and proprietary systems reveals marked performance gaps. Across all models, performance remains underwhelming on temporal reasoning, emotion recognition, and speech translation. Prompting in low-resource languages such as Burmese, Lao, Tamil, and Khmer lag behind English by over 5%. Our findings expose critical model limitations and underscore the need for inclusive model development. We will release datasets and the evaluation framework upon paper publication to facilitate reproducible benchmarking.

1 Introduction

Recent advancement in audio large language models (AudioLLMs) has led to transformative applications in voice assistants, transcription, accessibility technologies, and multimodal reasoning (Wu et al., 2024; Gemini Team, 2025; Zhang et al., 2023). Despite these advances, research in speech under-

standing has been disproportionately concentrated on high-resource languages, particularly English and a small number of European or East Asian languages (Yang et al., 2021; Wang et al., 2021; Bu et al., 2017). While recent benchmarking efforts (Sakshi et al., 2024; Wang et al., 2024; Yang et al., 2024) have made significant strides in evaluating audio-language models across diverse tasks and modalities, they universally overlook Southeast Asian (SEA) languages, leaving an entire linguistic region underexplored despite representing over 650 million speakers worldwide.

Developing comprehensive benchmarks for SEA languages also presents unique technical challenges. The region’s speech landscape is characterized by extraordinary linguistic diversity, rich tonal and phonetic structures, and substantial resource disparities across languages: factors that create evaluation complexities absent from English-centric benchmarks. Many SEA languages operate in low-resource contexts with limited annotated data and sparse digital representation, making robust evaluation both methodologically challenging and critically important for equitable technological development. While recent initiatives, such as MERaLiON (MERaLiON Team, 2024) which targets Singapore’s multilingual context, Typhoon2-Audio (Pipatanakul et al., 2024) which focuses on Thai, and SeaLLMs-Audio (Liu et al., 2025) which extends capabilities to selected SEA languages, have begun to build general-purpose speech-language models for the SEA region, these efforts remain limited in both scope and linguistic coverage. Crucially, they lack comprehensive evaluation frameworks necessary to systematically assess capabilities across the full spectrum of Southeast Asian speech understanding tasks. Meanwhile, fragmented data collection efforts across research groups have yielded heterogeneous SEA datasets (Lovenia and et al., 2024; Pham et al., 2023; Bustamin et al., 2024; Magic Data Technol-

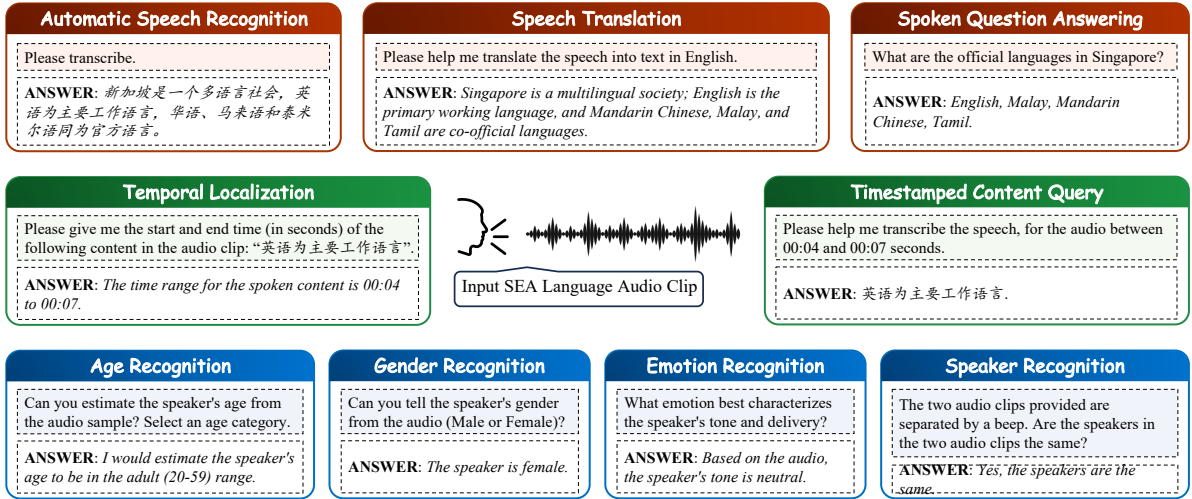


Figure 1: Overview of SEA-SPEECHBENCH task suite. The suite covers nine tasks in three categories. ■ **Speech processing**, ■ **Paralinguistic**, and ■ **Temporal understanding**.

ogy, 2025), but without unified frameworks for task definitions, prompting, and normalization, standardized comparison remains difficult.

In this work, we introduce the first ever comprehensive benchmark for speech understanding in Southeast Asian languages, designed to evaluate the capabilities of general-purpose speech-text LLMs. We focus on the official languages of Southeast Asian countries, as in Table 1, to balance broad regional coverage and practical relevance.

Table 1: Southeast Asian countries and their official language(s). Note: Several countries recognize additional regional or minority languages at sub-national levels; this table lists state-level official languages corresponding to the language codes used in our benchmark.

Country	Official language(s)	ISO 639-1 Code(s)
Singapore	English, Malay, Mandarin Chinese, Tamil	en, ms, zh, ta
Malaysia	Malay (Bahasa Melayu)	ms
Indonesia	Indonesian (Bahasa Indonesia)	id
Philippines	Filipino / Tagalog, English	tl, en
Thailand	Thai	th
Cambodia	Khmer (Cambodian)	km
Lao PDR	Lao	lo
Myanmar	Burmese (Myanmar)	my
Vietnam	Vietnamese	vi

Our benchmark encompasses 11 SEA languages with over 597 hours of audio data, curated from existing sources and synthesized into new tasks and datasets through systematic processing. The benchmark spans 9 diverse tasks across 3 broad categories: speech processing, paralinguistics, and temporal understanding. We introduce two temporal understanding tasks that evaluate models' capacity for temporal reasoning and localization within audio streams, addressing a previously unexplored dimension in audio LLM evaluation. These tasks test models' ability to navigate time-dependent in-

formation and extract content from specific temporal locations. The complete task suite is illustrated in Figure 1. To better reflect authentic usage scenarios, we evaluate each task using both English and native language text prompts.

The main contributions of this paper are threefold. First, we present SEA-SPEECHBENCH, the first ever large-scale multitask benchmark that systematically evaluates speech processing, paralinguistic analysis, and temporal understanding across Southeast Asian languages. It comprises a total of 99 evaluation sets with more than 97,000 audio samples. Second, we introduce temporal understanding tasks that assess models' ability to reason about time-dependent information in extended audio sequences. Finally, we comprehensively evaluate both open-source and proprietary models, offering critical insights into their strengths, limitations, and areas for further research.

2 SEA-SpeechBench: Task, Dataset and Evaluation Suite

2.1 Task Suite

SEA-SPEECHBENCH comprises 9 core tasks across 3 categories: *speech processing*, *paralinguistic analysis*, and *temporal understanding*. All tasks require models to produce textual responses given an audio input and a text query.

Speech processing covers three fundamental capabilities: Automatic Speech Recognition (ASR), Speech Translation (ST) from Southeast Asian languages to English, and Spoken Question Answering (SQA) based on SEA speech inputs.

Paralinguistic analysis examines vocal cues be-

Table 2: Summary of curated datasets for SEA-SPEECHBENCH. For multilingual datasets, each language-specific sub-dataset is counted separately in #Datasets.

Task	Languages	#Datasets	#Samples	Total L (h)	Min L (s)	Max L (s)
ASR	en tl id km lo ms my ta vi zh th	33	26,863	52.88	0.47	30.00
ST	en tl id km lo ms my vi th	9	7,189	26.46	3.06	30.00
SQA	en zh id th vi	7	5,462	40.57	20.00	30.00
ER	zh id th en ta	7	5,356	5.33	0.12	29.86
GR	zh id th en ta vi km my	16	13,599	22.02	0.12	29.90
AgeR	zh th en ta vi	5	4,608	6.55	0.58	20.78
SpkR	zh th en ta vi my	8	7,827	19.72	2.10	30.38
TCQ	zh en th id vi	7	13,145	211.98	20.00	180.00
TLoc	zh en th id vi	7	13,145	211.98	20.00	180.00
Total	–	99	97,194	597.49	–	–

yond linguistic content. It includes four tasks: Emotion Recognition (ER), which classifies emotional states; Gender Recognition (GR), which predicts gender from voice characteristics; Age Recognition (AgeR), which categorizes speakers as teens (10–19), adults (20–59), or seniors (60–100); and Speaker Recognition (SpkR), which determines whether two clips belong to the same speaker.

Temporal understanding introduces two novel tasks designed for extended audio, motivated by “skip to the content” and “what’s said at this time” use cases. Timestamped Content Query (TCQ) requires extracting content within a specified interval $[t_s, t_e]$, testing temporal grounding and localized retrieval. Temporal Localization (TLoc) asks models to predict the exact time span $\hat{y} = [\hat{t}_s, \hat{t}_e]$ where queried information appears, evaluating boundary detection and alignment. Together, these tasks formalize time-referenced retrieval and support fine-grained navigation across recordings.

The first two categories use short clips (≤ 30 s) to align with current model input limits. As real-world audio applications increasingly involve longer recordings where users require temporal navigation, SEA-SPEECHBENCH introduces the temporal understanding tasks to assess model ability to perform reasoning and localization over extended sequences of up to 3 minutes.

2.2 Data Curation

Dataset Statistics: As shown in Table 2, SEA-SPEECHBENCH is a comprehensive benchmark comprising of over 97,000 samples across 9 tasks

and 11 Southeast Asian languages.

Traditional speech processing tasks span the broadest linguistic coverage with standard 30-second clips, while our novel temporal understanding tasks focus on 5 key languages with extended recordings up to 180 seconds. Detailed information for each source dataset is provided in Appendix A. **Language and Task Distribution:** Figure 2 reveals the linguistic and task composition of our benchmark. The language distribution (Figure 2a) demonstrates substantial coverage across Southeast Asian languages, with English (18.9%), Thai (17.2%) and Chinese (16.8%) representing the largest segments, followed by Vietnamese (11.0%) and Indonesian (10.5%). We also include low-resource languages such as Khmer and Lao to ensure representation of the region’s full linguistic spectrum. Figure 2b shows balanced coverage across evaluation task families.

Audio Length Distribution: We cover both short and long audios to enable more comprehensive and realistic evaluation. As shown in Figure 2c, our benchmark encompasses a broad distribution of utterances spanning 0–30 seconds, with natural concentration around sentence-length segments for traditional tasks. For temporal understanding evaluation, we systematically curate extended recordings through strategic segmentation of long-form datasets, creating stratified duration bins of 30–60s, 60–120s, and 120–180s, totaling over 10,000 samples for long audios.

We applied uniform data processing across all datasets and synthesized additional evaluation sets

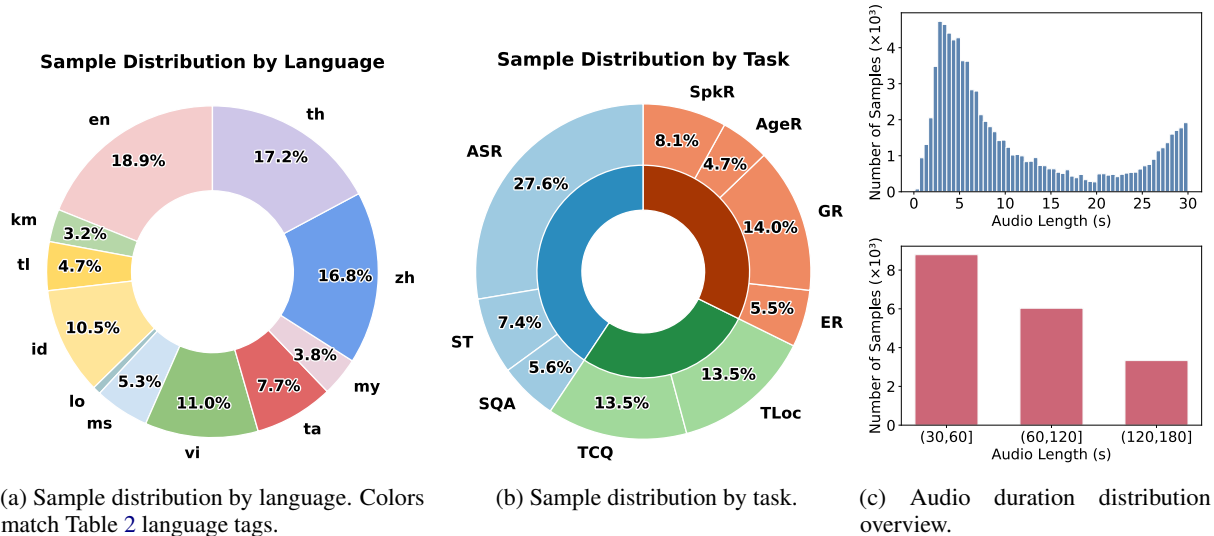


Figure 2: SEA-SPEECHBENCH composition: (a) language distribution by audio hours, (b) task distribution by sample count, and (c) audio duration distribution. (Upper: short audios (≤ 30 s) with fine-grained histogram. Lower: long audios (> 30 s) aggregated into three duration ranges.)

for tasks lacking suitable existing data. Processing and synthesis details are provided in Appendix B.

2.3 Metrics

Table 3 presents each task’s output format and primary evaluation metric. For ASR and TCQ, where the expected outputs are text transcripts, we employ script-dependent error rates: Character Error Rate (CER) for languages without explicit word boundaries (Chinese (zh), Thai (th), Khmer (km), Lao (lo), and Burmese (my)), and Word Error Rate (WER) for the other space-delimited languages. For ST, we use the BLEU score to measure n-gram overlap between generated translations and reference texts. For TLoc, given a actual time span $y = [t_s, t_e]$ and a prediction $\hat{y} = [\hat{t}_s, \hat{t}_e]$, define the intersection $I = [\min(\hat{t}_e, t_e) - \max(\hat{t}_s, t_s)]_+$ with $[x]_+ = \max(x, 0)$. Coverage and purity are $C = \frac{I}{t_e - t_s}$ and $P = \frac{I}{\hat{t}_e - \hat{t}_s}$, and the metric is $F1 = \frac{2CP}{C+P}$ for $C + P > 0$ (otherwise $F1 = 0$).

Table 3: Outputs and metrics by task. [†] *Judge-provided metrics*: scores are produced by a LLM serving as an external judge.

Task	Expected Output	Metric(s)
ASR	Text transcript	WER/CER
ST	Text translation	BLEU
SQA	Short textual answer	Scaled judge score s_{SQA} [†]
AgeR	Age bin $\in \{\text{teens, adults, seniors}\}$	Judge-based classification Acc [†]
ER	Emotion label	
GR	Gender label $\in \{\text{male, female}\}$	
SpkR	Speaker identity match or mismatch	
TCQ	Text transcript at given time	WER/CER
TLoc	Time span $[t_s, t_e]$	F1 score

Model-As-Judge: For the other tasks, due to the

inherent variability of free-form outputs, we employ an LLM-based judge for consistent evaluation, adapted based on the methodology proposed in (Wang et al., 2024). For SQA, the judge assigns a 0–5 quality score; we report this as a percentage via a linear scaling $s_{SQA} = 20 \times \text{score}$, $\text{score} \in \mathbb{Z}_{[0,5]}$. For paralinguistic tasks (ER, GR, AgeR, SpkR), the judge canonicalizes outputs and renders binary correctness, reported as accuracy in percent. Detailed judging prompts are provided in the Appendix F.

2.3.1 Assessment of Model Judge

To assess alignment with human judgment, we conduct an expert human evaluation for three model judges: Gemma3-27B-Instruct (Team, 2025), Qwen3-Omni-30B-A3B-Instruct (Xu et al., 2025b), and GPT-4.1 (OpenAI, 2024a), using a uniform evaluation protocol on a representative subset of tasks. We analyze judge–human agreement separately for English prompts and SEA language prompts. For each task and language, we randomly selected 100 model responses and collected the model judge scores and the human evaluator scores by native-speaker annotators. Alignment is quantified using Matthews correlation coefficient (MCC) for binary classification tasks and Spearman rank correlation (ρ) for ordinal ratings for SQA. Quantitative results are summarized in Table 4.

For all English prompts and binary classification settings, Gemma-3-27B Instruct and GPT-4.1 both exhibit strong and comparable alignment with human annotations. However, for SEA language prompts with ordinal rating scales, GPT-4.1

Table 4: Consistency between LLM-judged scores and human evaluator scores for English and SEA prompts. Alignment is quantified using Matthews correlation coefficient (MCC) for binary classification tasks and Spearman rank correlation (ρ) for ordinal ratings.

	MCC (Eng)	ρ (Eng)	MCC (SEA)	ρ (SEA)
Qwen3-Omni-30B-A3B-Instruct	0.81	0.80	0.76	0.37
Gemma-3-27B-Instruct	0.98	0.82	0.97	0.66
GPT-4.1	0.98	0.84	0.97	0.76

achieves substantially higher judge–human correlation than Gemma-3 ($\rho = 0.76$ vs. $\rho = 0.66$, averaged across SEA languages). Balancing both evaluation performance and computational cost, we therefore adopt a hybrid judging strategy: GPT-4.1 serves as the judge for SQA task, while Gemma-3-27B-Instruct is deployed in binary judge tasks.

3 Evaluations

3.1 Evaluation Setting

We conduct extensive evaluation across several state-of-the-art open-source audio/multimodal LLMs with multiple size variants ranging from 2B to 10B parameters: MERaLiON-2 (MERaLiON Team, 2024), SeaLLMs-audio (Liu et al., 2025), Phi-4-multimodal-instruct (Microsoft, 2025), Qwen2-Audio-Instruct (Chu et al., 2024), Qwen2.5-Omni (Xu et al., 2025a), Gemma-3n-it (Gemma Team, 2025), Voxtral (Voxtral Team, 2025), and Kimi-audio (KimiTeam, 2025). All models are evaluated using their official released checkpoints with recommended inference configurations to ensure fair comparison. To establish comprehensive performance baselines, we also evaluate two leading commercial models: Gemini 2.0 Flash (Gemini Team, 2025) and GPT-4o (OpenAI, 2024b). For GPT-4o, we employ the specialized Whisper-based transcription API for ASR tasks and the general audio understanding model for all other tasks, ensuring optimal performance.

3.2 Results and Insights

3.2.1 ASR Capability

We report the ASR performance of each model across different languages under English prompts in Table 7. Gemini 2.0 Flash has the best ASR performance overall, and MERaLiON-2 has the best performance amongst open-source models. We note that state-of-the-art performance on Tamil and Burmese is particularly low at more than 0.25 CER.

3.2.2 Speech Processing and Paralinguistics

Tables 5 and 6 present the results of additional speech processing and paralinguistic tasks. Results are reported under both English and SEA prompts. We observe performance deficiencies in ST where none of the open-source models achieved BLEU scores above 20. For paralinguistics tasks, emotion recognition is especially challenging. Apart from Kimi-Audio, no model achieved scores above 25. Even the best-performing Kimi-Audio did not score above 50.

Table 5: Results on speech processing tasks (ST, SQA) under English and SEA prompt respectively.

Model	Size	ST		SQA	
		ENG	SEA	ENG	SEA
Gemma-3n-it	2B	8.97	8.72	73.35	56.06
Qwen2.5-omni	3B	7.60	6.27	78.42	73.14
MERaLiON-2	3B	7.56	7.42	66.68	60.19
Voxtral	3B	19.98	18.15	82.32	79.06
Gemma-3n-it	4B	10.98	13.58	79.24	78.89
Phi-4	5.6B	3.04	0.32	64.69	47.36
SeaLLMs-Audio	7B	10.74	10.11	76.42	78.38
Qwen2-Audio-it	7B	4.54	3.20	65.82	57.95
Qwen2.5-omni	7B	7.91	8.06	71.60	75.57
Kimi-Audio	7B	3.36	7.71	69.62	63.49
MERaLiON-2	10B	17.75	19.52	82.00	82.05
Gemini 2.5 Flash	–	16.86	18.89	92.28	86.21
GPT-4o	–	21.24	21.39	86.66	82.44

Table 6: Results on paralinguistic tasks (AgeR, ER, GR, SpkR) under English and SEA prompt respectively.

Model	Size	AgeR		ER		GR		SpkR	
		ENG	SEA	ENG	SEA	ENG	SEA	ENG	SEA
Gemma-3n-it	2B	65.91	38.76	12.21	13.17	27.71	14.65	42.35	39.32
Qwen2.5-omni	3B	50.50	34.47	13.45	9.95	36.12	20.19	28.59	17.79
MERaLiON-2	3B	68.20	46.46	23.99	18.73	48.59	46.35	39.75	33.71
Voxtral	3B	75.71	52.61	10.62	5.35	27.65	9.89	42.90	36.09
Gemma-3n-it	4B	69.68	58.74	12.46	13.93	53.36	27.23	39.88	39.72
Phi-4	5.6B	43.45	41.98	20.87	9.20	44.21	32.60	23.27	27.48
SeaLLMs-Audio	7B	63.16	11.64	12.34	9.17	54.97	43.66	53.59	34.37
Qwen2-Audio-it	7B	20.21	23.13	24.47	19.36	92.25	66.31	49.85	36.56
Qwen2.5-omni	7B	20.15	23.13	16.33	10.45	50.46	34.58	14.63	11.44
Kimi-Audio	7B	47.91	48.42	36.86	42.55	90.97	79.23	54.16	48.16
MERaLiON-2	10B	66.24	57.31	18.73	20.34	59.93	46.28	53.25	46.12
Gemini 2.5 Flash	–	76.00	79.40	19.50	16.79	92.50	90.06	71.13	71.75
GPT-4o	–	76.60	61.20	17.00	19.50	46.13	31.63	10.50	8.38

3.2.3 Temporal Reasoning

In this section, we provide detailed analysis of temporal understanding capabilities, stratifying performance across four duration bins: [0,30), [30,60), [60,120), [120,180) seconds as shown in Table 8a. Figure 8b presents coverage (C), purity (P), and F_1 of TLoc task across audio-duration ranges, using MERaLiON-2-10B as a case study.

Metrics Analysis. First, we observe systematic over-coverage in temporal understanding across models. In TCQ, models frequently produce content that extends beyond the queried time window.

Table 7: **English-prompt ASR results averaged by language.** The performance is reported as WER or CER depending on language. Lower is better.

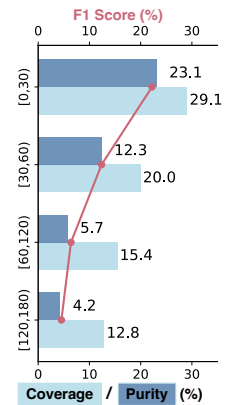
Model	Size	en	tl	vi	id	ta	th	zh	km	lo	ms	my
Gemma-3n-it	2B	1.32	0.42	3.77	0.17	1.20	3.03	1.55	3.92	1.42	1.06	2.81
Qwen2.5-omni	3B	0.07	0.60	0.27	0.15	1.36	0.12	0.05	2.97	1.39	0.23	2.01
MERaLiON-2	3B	0.07	0.20	0.35	0.11	0.43	0.08	0.11	1.72	0.86	0.17	1.19
Voxtral	3B	0.21	2.40	0.93	0.36	1.20	0.57	0.45	1.72	1.29	0.76	5.05
Gemma-3n-it	4B	0.76	0.24	2.71	0.51	0.45	0.32	0.69	1.71	0.15	1.98	2.98
SeaLLMs-Audio	7B	0.38	1.03	0.44	0.28	1.52	0.05	0.34	1.07	1.01	0.61	1.20
Phi-4	5.6B	0.09	5.10	2.92	2.76	1.93	2.98	0.10	3.11	2.45	3.15	2.13
Qwen2-Audio-it	7B	0.14	1.88	1.03	0.74	1.48	1.21	0.21	1.08	1.08	1.02	1.17
Qwen2.5-omni	7B	0.07	0.55	0.25	0.10	1.34	0.54	0.05	3.08	2.63	0.50	5.49
Kimi-Audio	7B	0.25	3.42	15.86	0.58	4.46	2.60	0.05	5.27	4.52	5.92	7.99
MERaLiON-2	10B	0.07	0.18	0.23	0.09	0.38	0.10	0.09	0.77	0.39	0.13	0.83
Gemini 2.5 Flash	-	0.11	0.13	0.12	0.03	0.27	0.05	0.18	0.13	0.16	0.11	0.31
GPT-4o	-	0.13	0.14	0.21	0.05	0.41	0.06	0.09	0.30	0.38	0.17	0.66

Table 8: Temporal understanding performance by duration.

(a) **English-prompt TCQ and TLoc results (%) by duration.** A dash (–) indicates audio lengths for which the model is unable to perform inference.

Model	Size	TCQ (WER/CER ↓)				TLoc (F1 Score ↑)			
		0–30 s	30–60 s	60–120 s	120–180 s	0–30 s	30–60 s	60–120 s	120–180 s
SeaLLMs-Audio	7B	5.48	-	-	-	11.57	-	-	-
Qwen2-Audio-it	7B	5.56	-	-	-	33.30	-	-	-
Qwen2.5-omni	3B	8.74	9.20	-	-	30.49	18.21	-	-
Gemma-3n-it	2B	7.32	7.49	-	-	11.82	8.27	-	-
Gemma-3n-it	4B	6.76	7.38	-	-	13.25	8.72	-	-
Phi-4	5.6B	20.67	14.90	16.25	-	12.97	6.25	3.64	-
MERaLiON-2	3B	4.77	8.21	11.27	-	18.82	10.28	5.14	-
Qwen2.5-omni	7B	5.49	6.58	9.85	-	35.74	19.98	11.32	-
Kimi-Audio	7B	14.00	19.07	24.04	42.78	14.49	8.61	3.70	3.00
Voxtral	3B	4.74	7.74	12.86	22.13	17.85	9.87	3.71	2.45
MERaLiON-2	10B	5.12	9.48	14.72	17.67	22.22	12.37	6.40	4.53
Gemini 2.5 Flash	-	2.41	2.64	9.44	4.09	11.66	7.70	5.57	5.30
GPT-4o	-	5.06	6.82	7.38	8.59	26.47	17.83	8.53	5.38

(b) TLoc metrics for MERaLiON-2-10B.



In TLoc, as demonstrated in Figure 8b, coverage consistently exceeds purity across all durations, which is a pattern we observe in most evaluated models. This asymmetry indicates weak temporal boundary localization and alignment, reflecting a recall-seeking strategy that favors longer spans, and thus higher coverage, at the cost of precision. These motivate finer-grained temporal grounding, boundary-aware training objectives, and decoding constraints that penalize span over-coverage.

Constraints on Audio Length. As shown in Table 8a, only a select subset of models sustains inference availability across all duration bins. Among open models, Voxtral, Kimi-Audio, and MERaLiON-2-10B demonstrate consistent availability, while commercial models handle the full range. Meanwhile, both TCQ and TLoc exhibit performance degradation with increasing duration: errors accumulate over longer contexts, manifesting as boundary drift and truncation, exposing current architectural limits in context window, frame compression, and long-range memory.

These findings underscore that temporal grounding represents an unresolved challenge in current audio-language architectures, with this deficiency becoming critically pronounced in long-form audio contexts where existing architectures prove inadequate for practical deployment.

3.2.4 Best-Performing Open-Source Models

To highlight model strengths across tasks and languages, we plot a winner map that marks, for each task-language pair, the best-performing **open-sourced** model in Figure 3. Each cell shows the top model and color-codes model identity. MERaLiON-2-10B establishes clear dominance in speech processing tasks, consistently achieving top performance across multiple languages including Indonesian, Vietnamese, and Filipino. Paralinguistic task leadership proves more distributed, with Kimi-Audio and Qwen2-Audio-Instruct alternately excelling in different linguistic contexts. Temporal understanding tasks reveal limited model coverage and inconsistent performance patterns, suggests

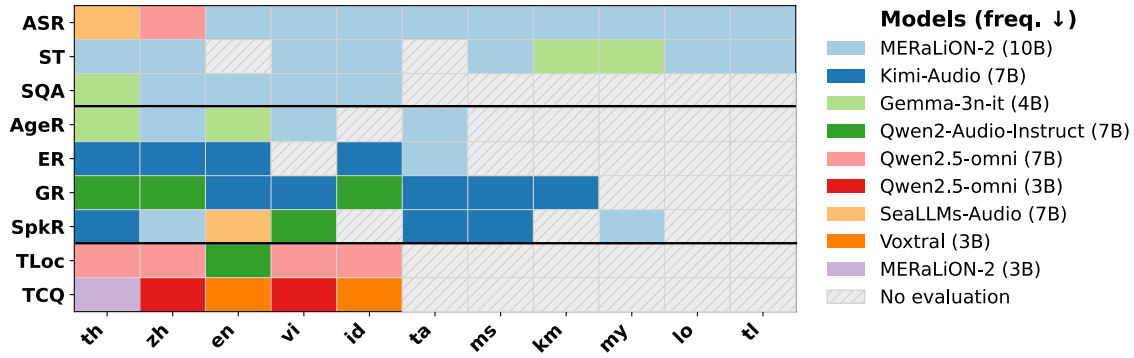


Figure 3: Winner map for open-sourced models across tasks and languages. Each cell marks the top model for a task–language pair; legend is ordered by overall win frequency.

that model strengths are specialized rather than generalizable, with no single model demonstrating comprehensive temporal reasoning proficiency.

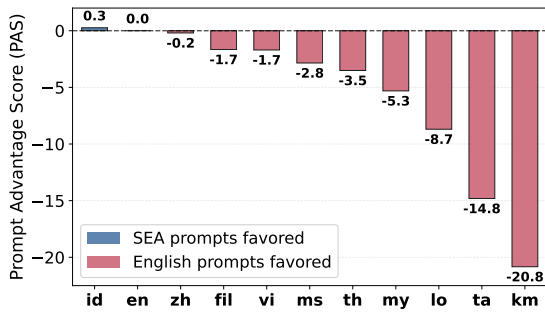


Figure 4: Cross-linguistic prompt sensitivity measured by Prompt Advantage Score (PAS).

3.2.5 Effect of Prompt Language

We systematically investigate how prompt language choice: English versus native SEA language, affects model performance across our benchmark.

Figure 4 demonstrates a cross-linguistic hierarchy in prompt sensitivity. We define a Prompt Advantage Score (PAS) to quantify this effect, with its detailed formulation provided in Appendix C. Higher PAS values indicate stronger local language prompt advantage, while negative scores suggest English prompt superiority for that particular language. Indonesian (id) emerges as the sole language showing consistent local prompt advantage (+0.3), while English (en) and Chinese (zh) exhibit near-neutral behavior (0.0, -0.2). The remaining SEA languages display increasing English preference across two distinct clusters: moderate disadvantages for Filipino, Vietnamese, Malay, Thai, and Myanmar (ranging from -1.7 to -5.3), and severe English advantages for Lao, Tamil, and Khmer (-8.7 to -20.8). This variation correlates strongly with orthographic and computational factors that

influence instruction parsing effectiveness:

Script Complexity and Tokenization. Non-Latin scripts create fundamental computational barriers. Languages like Thai, Lao, and Khmer lack clear word boundaries and employ complex grapheme clusters that disrupt standard tokenization processes. Abugida systems such as Myanmar and Tamil further complicate parsing through character-level ambiguities. These structural challenges impair instruction processing, while English prompts leverage well-established tokenization patterns that avoid such complications.

Training Data Quality and Orthographic Consistency. Corpus quality directly affects local prompt performance. Indonesian succeeds due to abundant Latin-script training data with consistent orthographic standards and imperative constructions, supporting reliable instruction following. In contrast, lower-resource languages suffer from limited, inconsistent training corpora marked by orthographic variations and dialectal diversity, making local prompts less reliable than standardized English alternatives.

These results expose a critical deployment gap: when users issue prompts in Southeast Asian languages, which is the natural interaction mode for regional populations, performance degrades substantially compared to English-prompted evaluation. This asymmetry underscores the urgent need for multilingual instruction-tuning that aligns with authentic user interaction patterns throughout model development and evaluation.

3.2.6 Refusal Behavior: Causes and Prevalence

From Table 6, GPT-4o attains only 10.50% accuracy on speaker recognition. Inspecting the error breakdown shows that these “errors” are refusals rather than wrong predictions. As illustrated in

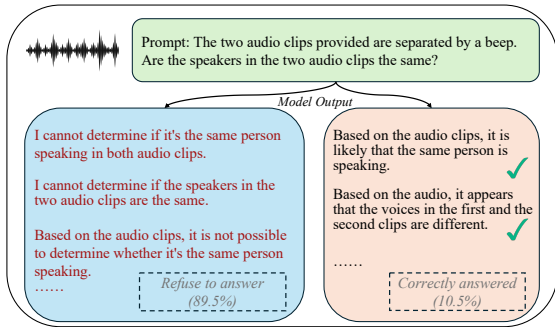


Figure 5: Speaker recognition failure examples in GPT-4o responses.

Figure 5, GPT-4o refuses to answer 89.5% of the queries. When the model does answer for 10.5% of the queries, it is consistently correct (non-refusal accuracy = 100%). This pattern emerges on our self-constructed SpkR dataset, which is likely out-of-distribution for GPT-4o, suggesting limited task generalization. Concurrently, the model may adopt a conservative, uncertainty-aware strategy rather than making overconfident predictions. By contrast, Qwen2.5-Omni-7B also attains low SpkR performance with frequent refusals, but its accepted responses contain nontrivial mistakes, pointing to weaker calibration and label grounding rather than abstention alone.

4 Related Works

Audio/Multimodal LLMs. Recent language models have evolved to process spoken audio alongside text through diverse architectural approaches. Early work established key paradigms: alignment-based models (Tang et al., 2024; Zhang et al., 2023; Wu et al., 2023) connect speech encoders to LLMs via lightweight adaptors, while unified decoders (Rubenstein et al., 2023; Nguyen et al., 2024) share token spaces for joint speech-text modeling. Contemporary systems demonstrate varied innovations: Phi-4 (Microsoft, 2025) employs mixture-of-LoRAs for multimodality, Qwen2.5-Omni (Xu et al., 2025a) advances temporal understanding through time-aligned position encoding, and Voxtral (Voxtral Team, 2025) handles extended recordings up to 40 minutes without separate ASR. Commercial systems like GPT-4o (OpenAI, 2024b) and Gemini (Gemini Team, 2025) provide end-to-end audio understanding within unified frameworks. For Southeast Asian languages, MERaLiON (MERaLiON Team, 2024) targets Singapore’s multilingual context while SeaLLMs-Audio (Liu et al., 2025) extends capabilities to five

major SEA languages. However, these efforts remain limited in scope, overlooking lower-resource languages and lacking comprehensive evaluation frameworks for the region’s full linguistic diversity.

Audio/Multimodal LLM Benchmarks. Audio-language evaluation has evolved from basic speech recognition to comprehensive multimodal assessment. Early efforts like SUPERB (Huang et al., 2024) aggregated speech tasks universally, and newer works (Wang et al., 2024; Yang et al., 2024) emphasize instruction-following across diverse audio types. Advanced benchmarks (Sakshi et al., 2024; Kumar et al., 2025) introduce complex reasoning with extended audio. Specialized evaluations target instruction-following (Gao et al., 2025) and domain-specific tasks (Ma et al., 2025). However, existing benchmarks remain predominantly English-centric with minimal Southeast Asian language coverage, creating a significant evaluation gap for low-resource linguistic contexts. This bias fundamentally limits understanding of model performance across the global linguistic landscape where inclusive deployment is most critically needed.

5 Conclusion

We present SEA-SPEECHBENCH, the first comprehensive benchmark for evaluating speech understanding across 11 Southeast Asian languages, comprising 97,000+ samples across 9 tasks in speech processing, paralinguistics, and temporal understanding. Our standardized framework enables reproducible, cross-linguistic comparisons through unified normalization and bilingual prompting.

Evaluation of leading commercial and open-source systems exposes systematic weaknesses: performance collapses on long audio (temporal brittleness), English prompts consistently outperform native languages (linguistic inequity), and tasks such as temporal reasoning, emotion recognition, and speech translation remain far below usability thresholds. These findings underscore persistent scalability and generalization gaps. By surfacing these limitations, SEA-SPEECHBENCH seeks to establish a rigorous baseline for developing temporally robust, linguistically inclusive, and practically deployable speech technologies for Southeast Asia’s diverse communities.

6 Limitations and Risks

Data availability remains a key limitation of this work. Despite extensive efforts to collect and curate data from diverse public and community-driven sources, coverage of certain dialects and task categories remains limited due to the fundamental scarcity of high-quality, annotated resources. This constraint restricts the breadth of dialectal and task-level evaluation that can be reliably supported. Future work may explore data-efficient learning paradigms, weakly supervised or synthetic data generation, and collaborative data collection initiatives to further extend coverage across underrepresented dialects and tasks.

Risk of test set contamination is another limitation. This risk is inherently difficult to control because most evaluated models are with undisclosed training data and curation pipelines.

References

- Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M. Tyers, and Gregor Weber. 2020. [Common voice: A massively-multilingual speech corpus](#). *Preprint*, arXiv:1912.06670.
- Hui Bu, Jiayu Du, Xingyu Na, Bengu Wu, and Hao Zheng. 2017. [Aishell-1: An open-source mandarin speech corpus and a speech recognition baseline](#). *Preprint*, arXiv:1709.05522.
- Anugrayani Bustamin, Andi M. Rizky, Elly Warni, Intan Sari Areni, and Indrabayu Indrabayu. 2024. [Indowavesentiment: Indonesian audio dataset for emotion classification](#). Mendeley Data, Version 1. CC BY 4.0 license; 300 audio files across five emotion categories.
- Ananlada Chotimongkol, Kwanchiva Saykhum, Patcharika Chootrakool, Nattanun Thatphithakkul, and Chai Wutiwiwatchai. 2009. [Lotus-bn: A thai broadcast news corpus and its research applications](#). In *2009 Oriental COCODA International Conference on Speech Database and Assessments*, pages 44–50.
- Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, Chang Zhou, and Jingren Zhou. 2024. [Qwen2-audio technical report](#). *arXiv preprint arXiv:2407.10759*.
- Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. 2022. [Fleurs: Few-shot learning evaluation of universal representations of speech](#). *arXiv preprint arXiv:2205.12446*.

- Yiming Gao, Bin Wang, Chengwei Wei, Shuo Sun, and AiTi Aw. 2025. [Ifeval-audio: Benchmarking instruction-following capability in audio-based large language models](#). *arXiv preprint arXiv:2505.16774*.
- Gemini Team. 2025. [Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities](#). *Preprint*, arXiv:2507.06261.
- Gemma Team. 2025. [Gemma 3 technical report](#). *Preprint*, arXiv:2504.18425.
- Fei He, Shan-Hui Cathy Chu, Oddur Kjartansson, Clara Rivera, Anna Katanova, Alexander Gutkin, Isin Demirsahin, Cibu Johny, Martin Jansche, Supheakmungkol Sarin, and Knot Pipatsrisawat. 2020. [Open-source Multi-speaker Speech Corpora for Building Gujarati, Kannada, Malayalam, Marathi, Tamil and Telugu Speech Synthesis Systems](#). In *Proceedings of The 12th Language Resources and Evaluation Conference (LREC)*, pages 6494–6503, Marseille, France. European Language Resources Association (ELRA).
- 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, and 1 others. 2024. [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)*, pages 12136–12140. IEEE.
- Yerbolat Khassanov, Zhiping Zeng, Van Tung Pham, Haihua Xu, and Eng Siong Chng. 2019. [Enriching rare word representations in neural language models by embedding matrix augmentation](#). In *Interspeech 2019*, page 3505–3509. ISCA.
- KimiTeam. 2025. [Kimi-audio technical report](#). *Preprint*, arXiv:2504.18425.
- Sonal Kumar, Šimon Sedláček, Vaibhavi Lokegaonkar, Fernando López, Wenyi Yu, Nishit Anand, Hyeonggon Ryu, Lichang Chen, Maxim Plička, Miroslav Hlaváček, and 1 others. 2025. [Mmau-pro: A challenging and comprehensive benchmark for holistic evaluation of audio general intelligence](#). *arXiv preprint arXiv:2508.13992*.
- Colin Leong, Joshua Nemecek, Jacob Mansdorfer, Anna Filighera, Abraham Owodunni, and Daniel Whitenack. 2022. [Bloom library: Multimodal datasets in 300+ languages for a variety of downstream tasks](#). *Preprint*, arXiv:2210.14712.
- Xinjian Li, Shinnosuke Takamichi, Takaaki Saeki, William Chen, Sayaka Shiota, and Shinji Watanabe. 2023. [Yodas: Youtube-oriented dataset for audio and speech](#). In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 1–8. IEEE.

617	Chaoqun Liu, Mahani Aljunied, Guizhen Chen,	OpenAI. 2024b. Gpt-4o system card . <i>Preprint</i> ,	671
618	Hou Pong Chan, Weiwen Xu, Yu Rong, and Wenxuan	arXiv:2410.21276.	672
619	Zhang. 2025. Seallms-audio: Large audio-language	Anh Pham, Khanh Linh Tran, Linh Nguyen, Thanh Duy	673
620	models for southeast asia. https://github.com/	Cao, Phuc Phan, and Duong A Nguyen. 2024.	674
621	DAM0-NLP-SG/SeaLLMs-Audio .	Bud500: A comprehensive vietnamese asr dataset .	675
622	Holy Lovenia and et al. 2024. Seacrowd: A multi-	Viet Thanh Pham, Xuan Thai Hoa Nguyen, Vu Hoang,	676
623	lingual multimodal data hub and benchmark suite	and Thi Thu Trang Nguyen. 2023. Vietnam-Celeb:	677
624	for southeast asian languages . <i>arXiv preprint arXiv:</i>	a large-scale dataset for Vietnamese speaker recog-	678
625	2406.10118 .	nition . In <i>Proc. INTERSPEECH 2023</i> , pages 1918–	679
626	Yinghao Ma, Siyou Li, Juntao Yu, Emmanouil Benetos,	1922.	680
627	and Akira Maezawa. 2025. Cmi-bench: A compre-	Kunat Pipatanakul, Potsawee Manakul, Natapong Ni-	681
628	hensive benchmark for evaluating music instruction	tarach, Warit Sirichotedumrong, Surapon Nonesung,	682
629	following. <i>arXiv preprint arXiv:2506.12285</i> .	Teetouch Jaknamon, Parinthapat Pengpun, Pittawat	683
630	Magic Data Technology. ASR-SgpCCSC .	Taveekitworachai, Adisai Na-Thalang, Sittipong Sri-	684
631	https://magichub.com/datasets/singaporean-	paisarnmongkol, Krisanapong Jirayoot, and Kasima	685
632	chinese-conversational-speech-corpus/ .	Tharnpipitchai. 2024. Typhoon 2: A family of	686
633	Magic Data Technology. 2023. Asr-sfdusc: A scripted	open text and multimodal thai large language models .	687
634	filipino daily-use speech corpus . MagicHub. 4.58	<i>Preprint</i> , arXiv:2412.13702.	688
635	hours of transcribed Filipino scripted speech focus-	Paul K. Rubenstein, Chulayuth Asawaroengchai,	689
636	ing on daily use sentences; 4,073 utterances by 10	Duc Dung Nguyen, Ankur Bapna, Zalán Borsos,	690
637	speakers. Licensed under CC BY-NC-ND 4.0.	Félix de Chaumont Quitry, Peter Chen, Dalia El	691
638	Magic Data Technology. 2025. Asr-malcs: Malay	Badawy, Wei Han, Eugene Kharitonov, Hannah	692
639	conversational speech corpus . Online: MagicHub.	Muckenhirn, Dirk Padfield, James Qin, Danny Rozen-	693
640	5 hours of transcribed Malay conversational speech,	berg, Tara Sainath, Johan Schalkwyk, Matt Shar-	694
641	10 conversations between 5 speaker pairs, licensed	ifi, Michelle Tadmor Ramanovich, Marco Tagliasac-	695
642	under CC BY-NC-ND 4.0.	chi, and 11 others. 2023. Audiopalm: A large lan-	696
643	MERaLiON Team. 2024. Meralion-audiollm: Bridg-	guage model that can speak and listen . <i>Preprint</i> ,	697
644	ing audio and language with large language models .	arXiv:2306.12925.	698
645	<i>Preprint</i> , arXiv:2412.09818.	S Sakshi, Utkarsh Tyagi, Sonal Kumar, Ashish Seth,	699
646	Microsoft. 2025. Phi-4-mini technical report: Com-	Ramaneswaran Selvakumar, Oriol Nieto, Ramani	700
647	pact yet powerful multimodal language models via	Duraiswami, Sreyan Ghosh, and Dinesh Manocha.	701
648	mixture-of-loras . <i>Preprint</i> , arXiv:2503.01743.	2024. Mmau: A massive multi-task audio under-	702
649	Myanmar Innovative Group. 2025. mig-burmese-	standing and reasoning benchmark . <i>arXiv preprint</i>	703
650	audio-transcription . https://huggingface.co/	<i>arXiv:2410.19168</i> .	704
651	datasets/Ko-Yin-Maung/mig-burmese-audio-	Keshan Sodimana, Knot Pipatsrisawat, Linne Ha, Mar-	705
652	transcription .	tin Jansche, Oddur Kjartansson, Pasindu De Silva,	706
653	Tu Anh Nguyen, Benjamin Muller, Bokai Yu,	and Supheakmungkol Sarin. 2018. A Step-by-Step	707
654	Marta R. Costa-jussa, Maha Elbayad, Sravya Pop-	Process for Building TTS Voices Using Open Source	708
655	uri, Christophe Ropers, Paul-Ambroise Duquenne,	Data and Framework for Bangla, Javanese, Khmer,	709
656	Robin Algayres, Ruslan Mavlyutov, Itai Gat, Mary	Nepali, Sinhala, and Sundanese . In <i>Proc. The 6th</i>	710
657	Williamson, Gabriel Synnaeve, Juan Pino, Benoit	<i>Intl. Workshop on Spoken Language Technologies for</i>	711
658	Sagot, and Emmanuel Dupoux. 2024. Spirit lm:	<i>Under-Resourced Languages (SLTU)</i> , pages 66–70,	712
659	Interleaved spoken and written language model .	Gurugram, India.	713
660	<i>Preprint</i> , arXiv:2402.05755.	Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao	714
661	Yin May Oo, Theeraphol Wattanavekin, Chenfang Li,	Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, and	715
662	Pasindu De Silva, Supheakmungkol Sarin, Knot Pi-	Chao Zhang. 2024. Salmonn: Towards generic hear-	716
663	patsrisawat, Martin Jansche, Oddur Kjartansson, and	ing abilities for large language models . <i>Preprint</i> ,	717
664	Alexander Gutkin. 2020. Burmese Speech Corpus,	arXiv:2310.13289.	718
665	Finite-State Text Normalization and Pronunciation	Gemma Team. 2025. Gemma 3 technical report .	719
666	Grammars with an Application to Text-to-Speech . In	<i>Preprint</i> , arXiv:2503.19786.	720
667	<i>Proceedings of The 12th Language Resources and</i>	Thanushs25. 2024. tamil-audio-emotion-	721
668	<i>Evaluation Conference (LREC)</i> .	classification . https://huggingface.co/	722
669	OpenAI. 2024a. Gpt-4 technical report . <i>Preprint</i> ,	datasets/Thanushs25/tamil-audio-emotion-	723
670	arXiv:2303.08774.	classification .	724
		Jubeerathan Thevakumar, Luxshan Thavarasa, Thanikan	725
		Sivatheepan, Sajeew Kugarajah, and Uthayasanker	726

727	Thayasivam. 2025. Emota: A tamil emotional speech dataset . In <i>Proceedings of the First Workshop on Challenges in Processing South Asian Languages (CHI PSAL 2025)</i> , pages 193–201, Abu Dhabi, UAE. International Committee on Computational Linguistics.		
728			
729			
730			
731			
732			
733	Voxtral Team. 2025. Voxtral . <i>Preprint</i> , arXiv:2507.13264.		
734			
735	Hoang Long Vu, Phuong Tuan Dat, Pham Thao Nhi, Nguyen Song Hao, and Nguyen Thi Thu Trang. 2025. Voxvietnam: a large-scale multi-genre dataset for vietnamese speaker recognition . In <i>ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 1–5.		
736			
737			
738			
739			
740			
741	Bin Wang, Xunlong Zou, Geyu Lin, Shuo Sun, Zhuohan Liu, Wenyu Zhang, Zhengyuan Liu, AiTi Aw, and Nancy F Chen. 2024. Audiobench: A universal benchmark for audio large language models. <i>arXiv preprint arXiv:2406.16020</i> .		
742			
743			
744			
745			
746	Changhan Wang, Morgane Rivière, Ann Lee, Anne Wu, Chaitanya Talnikar, Daniel Haziza, Mary Williamson, Juan Pino, and Emmanuel Dupoux. 2021. Voxpopuli: A large-scale multilingual speech corpus for representation learning, semi-supervised learning and interpretation . <i>Preprint</i> , arXiv:2101.00390.		
747			
748			
749			
750			
751			
752	Jian Wu, Yashesh Gaur, Zhuo Chen, Long Zhou, Yimeng Zhu, Tianrui Wang, Jinyu Li, Shujie Liu, Bo Ren, Linqun Liu, and Yu Wu. 2023. On decoder-only architecture for speech-to-text and large language model integration . <i>Preprint</i> , arXiv:2307.03917.		
753			
754			
755			
756			
757			
758	Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. 2024. Next-gpt: Any-to-any multimodal llm . <i>Preprint</i> , arXiv:2309.05519.		
759			
760			
761	Li Xinjian, Takamichi Shinnosuke, Saeki Takaaki, Chen William, Shiota Sayaka, and Watanabe Shinji. 2023. Yodas: Youtube-oriented dataset for audio and speech . In <i>2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)</i> , pages 1–8. IEEE.		
762			
763			
764			
765			
766			
767	Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan, Kai Dang, Bin Zhang, Xiong Wang, Yunfei Chu, and Junyang Lin. 2025a. Qwen2.5-omni technical report . <i>arXiv preprint arXiv:2503.20215</i> .		
768			
769			
770			
771			
772	Jin Xu, Zhifang Guo, Hangrui Hu, Yunfei Chu, Xiong Wang, Jinzheng He, Yuxuan Wang, Xian Shi, Ting He, Xinfu Zhu, Yuanjun Lv, Yongqi Wang, Dake Guo, He Wang, Linhan Ma, Pei Zhang, Xinyu Zhang, Hongkun Hao, Zishan Guo, and 19 others. 2025b. Qwen3-omni technical report . <i>Preprint</i> , arXiv:2509.17765.		
773			
774			
775			
776			
777			
778			
779	Qian Yang, Jin Xu, Wenrui Liu, Yunfei Chu, Ziyue Jiang, Xiaohuan Zhou, Yichong Leng, Yuanjun Lv,		
780			
		Zhou Zhao, Chang Zhou, and 1 others. 2024. Airbench: Benchmarking large audio-language models via generative comprehension . <i>arXiv preprint arXiv:2402.07729</i> .	781 782 783 784
		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, and 1 others. 2021. Superb: Speech processing universal performance benchmark . <i>Interspeech 2021</i> .	785 786 787 788 789 790
		Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. 2023. Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities . <i>Preprint</i> , arXiv:2305.11000.	791 792 793 794 795
		Jinming Zhao, Tengan Zhang, Jingwen Hu, Yuchen Liu, Qin Jin, Xinchao Wang, and Haizhou Li. 2022. M3ed: Multi-modal multi-scene multi-label emotional dialogue database . <i>arXiv preprint arXiv:2205.10237</i> .	796 797 798 799 800
		Kun Zhou, Berrak Sisman, Rui Liu, and Haizhou Li. 2022. Emotional voice conversion: Theory, databases and esd . <i>Speech Communication</i> , 137:1–18.	801 802 803 804
		A Dataset Catalog	805
		Table 9 lists the source datasets used in our SEA-SpeechBench, along with their tasks, languages, audio lengths and licensing information.	806 807 808
		B Data Processing and Synthesis	809
		We provide additional details on data processing and the construction of new datasets from source materials.	810 811 812
		B.1 Data Processing	813
		Audio and Prompt Standardization: All audio was resampled to 16 kHz for consistency and model compatibility. We provide parallel prompts in native SEA languages and English to reflect realistic usage patterns, constructing standardized templates where needed. Prompt examples are provided in the Appendix E.	814 815 816 817 818 819 820
		Dataset Sampling: 1. We adopted official test splits when available or sampled 1,000 instances while preserving identifiers to prevent contamination. 2. For our evaluation set, we further sampled 1,000 instances per dataset using class-balanced sampling for classification tasks, ensuring computational feasibility while maintaining representativeness. 3. We checked that uniquely identifying information such as names were removed when	821 822 823 824 825 826 827 828 829

Table 9: Overview of datasets: Dataset statistics, including language coverage, task types, basic audio properties, and licensing information. *Including long audios (>30s).

	Languages	Tasks	Total L (hr)	Min L (s)	Max L (s)	License
Commonvoice (Ardila et al., 2020)	zh, vi, th, id, ta, en	GR, AGE, ASR	22.66	0.58	20.78	MPL 2.0
FLEURS (Conneau et al., 2022)	zh, vi, th, my, ms, lo, km, id, tl	ASR, ST, GR	57.79	3.06	30.00	CC-BY 4.0
OpenSLR (Sodimana et al., 2018) (He et al., 2020; Oo et al., 2020)	ta, my, km	GR, ASR	6.33	1.71	17.15	CC-BY-SA 4.0
Bloom-Speech (Leong et al., 2022)	tl, my, en	ASR	0.84	0.47	29.46	CC BY-NC 4.0
Thai Elderly Speech (Lovenia and et al., 2024)	th	ASR, GR, SpkR	5.66	2.04	27.06	CC-BY-SA 4.0
LOTUS (Chotimongkol et al., 2009)	th	ASR	1.22	1.84	29.78	CC-BY-NC-SA 4.0
THAI SER (Lovenia and et al., 2024)	th	SpkR, GR, ER	6.22	0.60	29.86	CC-BY-SA 4.0
SG Streets (Khassanov et al., 2019)	en	TLoc, TCQ, SQA, ASR, GR	6.55	0.80	59.97	NA
VietMed (Pham et al., 2023)	vi	ASR	1.74	2.00	12.00	MIT
VoxVietnam-O (Vu et al., 2025)	vi	SpkR	2.43	1.42	29.68	CC
Bud500 (Pham et al., 2024)	vi	ASR	0.71	1.01	5.48	Apache-2.0
Vietnam-Celeb (Pham et al., 2023)	vi	GR	2.14	0.84	29.66	CC-BY-4.0
ASR-SMaDuSC (Lovenia and et al., 2024)	ms	ASR, SpkR, GR	8.51	2.64	29.22	CC BY-NC-ND 4.0
ASR-MALCSC (Magic Data Technology, 2025)	ms	ASR	0.74	0.61	18.16	CC BY-NC-ND 4.0
IndoWaveSentiment (Bustamin et al., 2024)	id	GR, ER	0.52	3.00	3.80	CC-BY-4.0
EmoTa (Thevakumar et al., 2025)	ta	SpkR, GR, ER	2.90	1.06	9.79	EACL
SFDUSC (Magic Data Technology, 2023)	tl	SpkR, GR, ASR	4.86	2.03	21.49	CC BY-NC-ND 4.0
YODAS2* (Li et al., 2023)	en, id, zh, th, vi	TLoc, TCQ, SQA	493.12	20.01	180.00	CC-BY 3.0
ESD (Zhou et al., 2022)	en, zh	ASR, ER, SpkR	6.56	1.41	10.79	MIT
M3ED (Zhao et al., 2022)	zh	ER, GR	0.83	0.12	7.04	CC BY-NC-ND 4.0
MIG (Myanmar Innovative Group, 2025)	my	ASR, SpkR	1.50	0.80	21.39	NA
TEC (Thanush25, 2024)	ta	ER	0.68	2.69	29.86	NA
ASR-SgpCCSC (Magic Data Technology)	zh	ASR	43.76	0.88	180.00	CC BY-NC-ND 4.0

curating and processing data from open-source materials. We also applied profanity and content filtering on raw speech transcripts in the YODAS2 dataset.

This standardized processing transforms diverse datasets into a unified, reproducible benchmark supporting consistent evaluation across models.

B.2 New Evaluation Data Construction

SpkR (Speaker Recognition): We constructed datasets by sampling and concatenating pairs of audio segments, separated by a short beep, drawn either from the same speaker or from different speakers, to create controlled instances for the speaker verification task. Refer to Table 9 for the datasets used to synthesize the speaker recognition task.

SQA (Spoken Question Answering): We carefully filtered the YODAS2 dataset (Xinjian et al., 2023) using CTC forced alignment scores between the audio clips and the provided transcriptions with a log-probability threshold of -1, exact language match between the audio and text transcription, as well as profanity and content filtering. Based on the cleaned transcriptions, we generated question-answer pairs using GPT-4.1 (OpenAI, 2024b).

TCQ (Timestamped Content Query): We filtered the YODAS2 dataset following the same procedure used for SQA. We then identified utterances whose CTC log-probability exceeded -0.1, recording their start and end timestamps. The task was formulated such that models are required to tran-

scribe utterances occurring between the specified timestamps.

TLoc (Temporal Localization): We adopted the same preprocessing procedure as in TCQ. The task was formulated such that models are required to predict the start and end timestamps corresponding to specified utterances.

In all cases, we applied the same standardized audio preprocessing, sample selection, and prompt construction pipeline, ensuring these synthesized datasets are consistent and compatible with the broader benchmark.

C Definition of Prompt Advantage Score (PAS)

We quantify prompt-language effects with a scale-invariant score capturing both *magnitude* and *direction*. Let s_{SEA} and s_{EN} denote task scores for SEA language prompt and English prompt, respectively. For each model M , task T , and language L , define the symmetric relative difference

$$d_{M,T,L} = \frac{|s_{SEA} - s_{EN}|}{\frac{|s_{SEA}| + |s_{EN}|}{2} + \epsilon}, \quad (1)$$

and aggregate across models within task by the median $\tilde{d}_{T,L} = \text{median}_M d_{M,T,L}$. Directionality is encoded by the model-wise win rate of local prompts $w_{T,L} = \mathbb{P}_M(s_{SEA} > s_{EN})$. The task-level Prompt Advantage Score combines direction and

886 effect size:

$$887 \quad \text{PAS}_{T,L} = (2(w_{T,L} - 0.5)) \cdot \tilde{d}_{T,L}, \quad (2)$$

888 so $\text{PAS}_{T,L} > 0$ favors local prompts and
889 $\text{PAS}_{T,L} < 0$ favors English. The language-level
890 score averages over the tasks covered by L , denoted
891 $\mathcal{T}(L)$:

$$892 \quad \text{PAS}_L = \frac{1}{|\mathcal{T}(L)|} \sum_{T \in \mathcal{T}(L)} \text{PAS}_{T,L}. \quad (3)$$

893 By construction, $|\text{PAS}_L|$ measures the strength of
894 prompt sensitivity while $\text{sign}(\text{PAS}_L)$ indicates the
895 preferred prompt language.

896 D Postprocessing of Evaluation

897 To ensure fair and consistent evaluation across
898 Southeast Asian languages, we design a unified text
899 normalization pipeline comprising three sequential
900 stages:

901 **Canonical Form Conversion.** Transcripts undergo
902 initial standardization through NFC Unicode nor-
903 malization and systematic lowercasing to eliminate
904 encoding inconsistencies and case variations. Con-
905 figurable punctuation filtering preserves linguisti-
906 cally meaningful characters (hyphens, apostrophes)
907 while removing extraneous symbols that could in-
908 troduce evaluation noise.

909 **Content Filtering and Lexical Standardization.**
910 Non-linguistic artifacts including speaker tags,
911 bracketed annotations, and conversational fillers
912 are systematically removed. Digital content under-
913 goes targeted normalization, and English contrac-
914 tions receive consistent expansion. For character-
915 based writing systems, e.g., Chinese and Thai, inter-
916 character spacing is inserted to enable reliable tok-
917 enization.

918 **Language-Specific Enhancement and Validation.**
919 We collaborated with native speakers across all
920 SEA languages to validate and refine normalization
921 heuristics.

922 E Evaluation Prompt Examples

923 We evaluate using parallel prompts in native SEA
924 languages and English. Prompt examples are pro-
925 vided in Figure 6 and 7.

Figure 6: Overview of prompt examples: illustrative prompts provided for ASR, ST, SQA, AgeR, and ER tasks across multiple languages.

Tasks	Prompt Language	Input Example	Answer Example
Automatic Speech Recognition	ENG	Text: Please help me transcribe the speech into text in English. 🗣️: [English Speech]	🗣️: Activity. Draw a picture of you playing a sport.
		Text: Listen to the speech and provide the text version in Malay. 🗣️: [Malay Speech]	🗣️: Angin kencang, hujan batu, pemendakan berlebihan dan kebakaran hutan adalah bentuk dan kesan cuaca yang teruk. Begitu juga ribut petir, tornado, semburan air dan siklon.
	SEA	Text: Sila bantu saya transkripsikan ucapan ini ke dalam bentuk teks dalam Bahasa Melayu. 🗣️: [Malay Speech]	🗣️: Angin kencang, hujan batu, pemendakan berlebihan dan kebakaran hutan adalah bentuk dan kesan cuaca yang teruk. Begitu juga ribut petir, tornado, semburan air dan siklon.
		Text: กรุณาช่วยถอดเสียงพูดนี้เป็นข้อความในภาษาไทยด้วย 🗣️: [Thai Speech]	🗣️: มันเกิดอุบัติเหตุที่ท่ารถ ทำให้รถเกือบชนคนเสียประจวบ
Speech Translation	ENG	Text: Please help me translate the speech into text in English 🗣️: [Malay Speech]	🗣️: The television reports that the white smoke seen is from plants.
		Text: Recognize the verbal content in the speech and translate it into text in English. 🗣️: [Thai Speech]	🗣️: The smaller the Rossby number, the less movement the star will have in relation to the reversal of the Earth's magnetic poles.
	SEA	Text: အင်္ဂလိပ် သို့ စကားပြောချက်ကို စာသားအဖြစ် ဘာသာပြန်ပေးရန် ကျေးဇူးပြုပါ။ 🗣️: [Burmese Speech]	🗣️: Strong winds, heavy rain and sleet, and lightning are forms and effects of severe weather such as thunderstorms, tornadoes, waterspouts, and cyclones.
		Text: Vui lòng giúp tôi chuyển lời nói thành văn bản bằng Tiếng Anh 🗣️: [Vietnamese Speech]	🗣️: In his writing on the presidential speech, Oliver Sacks points out how people who are unable to understand the speech due to brain damage can still accurately judge its truthfulness.
Spoken Question Answering	ENG	Text: How many housing units were scheduled to be built in Bedok under the mentioned plan? 🗣️: [English Speech]	🗣️: Sixteen thousand five hundred housing units.
		Text: What can the playback tab be used for? 🗣️: [English Speech]	🗣️: It can be used to play videos or adjust effects.
	SEA	Text: Untuk apa sendok diambil dalam proses mengambil air sirih merah? 🗣️: [Indonesian Speech]	🗣️: Sendok diambil agar bisa merendam semua sirih merah.
		Text: Vì sao người nói cảm thấy buồn đau và phải chôn giấu nỗi buồn đó? 🗣️: [Vietnamese Speech]	🗣️: Vì đi qua nhiều mối tình dở vỡ mà không ai ở lại, khiến người nói cảm thấy buồn đau và phải chôn giấu nỗi buồn đó.
Age Recognition	ENG	Text: Can you guess the speaker's age based purely on their voice characteristics? Choose among these age categories: teens (10-19), adult (20-59), or senior (60-100). 🗣️: [English Speech]	🗣️: adult (20-59)
		Text: Can you estimate the speaker's age from the audio sample? Select an age category: teens (10-19), adult (20-59), or senior (60-100). 🗣️: [Thai Speech]	🗣️: Based on the audio sample, the speaker is likely in the adult (20-50) age category.
	SEA	Text: பேச்சை அடிப்படையாகக் கொண்டு, பேசும் நபருக்கு ஏற்ற சராசரி வயது வரம்பை வழங்க முடியுமா? தயவுசெய்து ஒரு வயது பிரிவை தேர்ந்தெடுக்கவும்: இளையோர் (10-19), பெரியோர் (20-59), அல்லது மூத்தோர் (60-100). 🗣️: [Tamil Speech]	🗣️: பெரியோர் (20-59)
		Text: สามารถเดาอายุของผู้พูด จากเสียงพูดเพียงอย่างเดียวหรือกลุ่มอายุ? กรุณาเลือกหนึ่งกลุ่มอายุ จาก: วัยรุ่น (10-19), ผู้ใหญ่ (20-59), ผู้สูงอายุ (60-100) 🗣️: [Thai Speech]	🗣️: ผู้ใหญ่ (20-59)
Emotion Recognition	ENG	Text: Based on the speaker's speech patterns, what do you think they are feeling? 🗣️: [English Speech]	🗣️: The speaker's speech suggesting they might be feeling neutral.
		Text: Can you identify or describe any emotions or feelings expressed by the speaker? Answer only using one sentence. Do not explain. 🗣️: [Tamil Speech]	🗣️: happy
	SEA	Text: คุณสามารถระบุอารมณ์ของอารมณ์หรือความรู้สึกที่ผู้พูดแสดงออกมาได้หรือไม่? กรุณาตอบด้วยประโยคเดียวเท่านั้น ห้ามอธิบายเพิ่มเติม 🗣️: [Thai Speech]	🗣️: ผู้พูดแสดงความรู้สึกหงุดหงิด ไม่พอใจ และอารมณ์ขุ่นมัวหรืออาจรุนแรง
		Text: Apa interpretasi Anda terhadap emosi yang ditunjukkan dari petunjuk emosional yang ada dalam audio? Berikan jawaban satu kalimat tentang emosi tersebut. Jangan berikan alasan. 🗣️: [Indonesian Speech]	🗣️: Emosi yang ditunjukkan adalah kepuasan.

Figure 7: Overview of prompt examples: illustrative prompts provided for GR, SpkR, TCQ, and TLoc tasks across multiple languages.

Tasks	Prompt Language	Input Example	Answer Example
Gender Recognition	ENG	Text: Can you identify the speaker's gender based on the audio (Male or Female)? 🗣️: [English Speech]	🗣️: male.
		Text: From the audio, can you guess the speaker's gender (Male or Female)? 🗣️: [Thai Speech]	🗣️: female.
	SEA	Text: ஆடியோவை அடிப்படையாகக் கொண்டு பேசும் நபரின் பாலினத்தை (ஆண் அல்லது பெண்) அடையாளம் காண முடியுமா? 🗣️: [Tamil Speech]	🗣️: ஆண்.
		Text: Dapatkah Anda membedakan jenis kelamin pembicara berdasarkan audio ini (Laki-laki atau Perempuan)? 🗣️: [Indonesian Speech]	🗣️: Laki-laki.
Speaker Recognition	ENG	Text: A beep sound is placed between the two audio segments. Is there a match between the identity of the speakers in the two recordings? 🗣️: [English Speech]	🗣️: Yes, the speakers are the same.
		Text: A short beep marks the break between the two audio clips. Do you think both clips feature the same speaker? 🗣️: [Burmese Speech]	🗣️: No, the speakers are different.
	SEA	Text: Bunyi beep diletakkan di antara dua segmen audio. Adakah terdapat padanan identiti penutur dalam kedua-dua rakaman ini? 🗣️: [Malay Speech]	🗣️: Tidak, pembesar suara berbeza.
		Text: Hai đoạn âm thanh được cung cấp được ngăn cách bởi một tiếng bíp. Liệu người nói trong hai đoạn âm thanh này có phải là cùng một người không? 🗣️: [Vietnamese Speech]	🗣️: Không, loa thì khác nhau.
TCQ	ENG	Text: Please help me transcribe the speech into text in English, for the audio between 21.96 and 25.54 seconds. 🗣️: [English Speech]	🗣️: The first phase of the eastern line which included aljunied.
		Text: Process the speech and provide the text output in Thai for the duration 6.07 to 9.6 seconds. 🗣️: [Thai Speech]	🗣️: การซ่อมที่สี่ทุ่มสองซ่อมเสร็จแล้วกว่าสี่ทุ่ม
	SEA	Text: Nghe đoạn phát biểu và cung cấp phiên bản văn bản bằng Tiếng Việt, bao gồm phần từ 12.03 đến 16.38 giây. 🗣️: [Vietnamese Speech]	🗣️: Tinh thuế cho từng cá nhân hoàn toàn tự động để nộp về từng quốc gia sẽ tại.
		Text: Ubah masukan ucapan menjadi transkripsi teks dalam Bahasa Indonesia. Bagian audio yang relevan adalah dari detik 15.12 hingga 20.87. 🗣️: [Indonesian Speech]	🗣️: Kemudian untuk voice setting di sini ada lagi di sini kalian bisa menambahkan untuk kecepatannya.
TLoc	ENG	Text: Please give me the start and end time (in seconds) of the following content in the audio clip: <i>80% equal to 1.38 kilograms.</i> 🗣️: [English Speech]	🗣️: The time interval is **[00:22:257 - 00:25:437]** .
		Text: Identify when the following content occurs in the audio by providing the start and end time in seconds: <i>dan sekian informasi dari saya dan kiranya bilamana ada yang ditanyakan silahkan.</i> 🗣️: [Indonesian Speech]	🗣️: The start and end time of the content is: \n**Start:** 00:17.060\n**End:** 00:22.259
	SEA	Text: : ช่วงเวลาที่ตรงกับเนื้อหาในไฟล์เสียง (เป็นวินาที): การออกแบบเสื้อผ้าเครื่องแต่งกาย. 🗣️: [Thai Speech]	🗣️: [0:0:165 – 0:1:835] การออกแบบเสื้อผ้าเครื่องแต่งกาย
		Text: Kapan frasa ini muncul di audio? Berikan waktu mulai dan selesai dalam detik: dan sekian informasi dari saya dan kiranya bilamana ada yang ditanyakan silahkan. 🗣️: [Indonesian Speech]	🗣️: Tentu, frasa tersebut muncul pada waktu berikut:\nMulai: 17:039 detik\nSelesai: 22:189 detik

F Model-as-Judge Prompt Examples

Table 10: Model-as-Judge prompts, adopted from (Wang et al., 2024). To ensure protocol consistency, human judges followed the same evaluation instruction.

Task	Prompt
SQA	<p>[Question] {question}</p> <p>[Reference] {reference}</p> <p>[Model Prediction] {prediction}</p> <p>[Task] Rate the model prediction based on its alignment with the reference, focusing on accuracy and relevance to the reference. Be critical. Score0: The prediction repeats or rephrases the question without giving an answer. Score0: The prediction is refusing to give concrete results, providing something like 'cannot decide'. Score0: The prediction is completely misaligned, providing incorrect or irrelevant information compared to the reference. Score1: The prediction shows minimal alignment, often misunderstanding or providing irrelevant details unrelated to the reference. Score2: The prediction recognizes the topic but diverges significantly from the reference in accuracy or relevance. Score3: The prediction aligns with the reference generally but lacks detail or precise accuracy in some aspects. Score4: The prediction is mostly accurate and relevant, closely following the reference but could be clearer or more detailed. Score5: The prediction is highly accurate, detailed, and matches the reference perfectly, capturing its essence and detail.</p> <p>Your response should be formatted as follows: Explanation: (Provide a concise explanation of your rating, comparing the reference with the model prediction. "The reference is [XXX], while the model prediction is [YYY]. I think ...") Rating: (int)</p>
AgeR ER GR SpkR	<p>[Question] {question}</p> <p>[Reference] {reference}</p> <p>[Model Prediction] {prediction}</p> <p>[Task] Rate the model prediction based on its alignment with the reference, focusing on accuracy and relevance to the reference. Be critical. Score0: The prediction repeats or rephrases the question without giving an answer. Score0: The prediction is refusing to give concrete results, providing something like 'cannot decide'. Score0: The prediction is wrong, providing incorrect or irrelevant information compared to the reference. Score1: The prediction is correct, capturing or covering the meaning from the reference.</p> <p>Your response should be formatted as follows: Explanation: (Provide a concise explanation of your rating, comparing the reference with the model prediction. "The reference is [XXX], while the model prediction is [YYY]. I think ...") Rating: (int)</p>

Table 11: ASR performance: results across datasets comparing English and SEA prompts. Scores are reported as raw WER and CER values without normalization; lower values indicate better performance.

Model				MERaLiON2 10B	MERaLiON2 3B	SeaLLMs Audio 7B	Phi-4 multi- modal instruct	Kimi Audio	Voxtral mini	Qwen2 Audio 7B Instruct	Qwen 2.5 Omni 3B	Qwen 2.5 Omni 7B	gemma 3n E4B-it	gemma 3n E2B-it	Gemini 2.5 Flash	Whisper large v3	GPT 4o Audio	
size				10B	3B	7B	5.6B	7B	3B	7B	3B	7B	4B	2B	-	1.5B	-	
Data	Metrics	Lang	Prompt	No Prompt														
Bloom-Speech	WER	en	ENG	0.06	0.07	0.74	0.06	0.06	0.20	0.13	0.06	0.06	0.38	0.52	0.05	0.04	0.06	
			SEA	0.06	0.07	0.74	0.06	0.06	0.20	0.13	0.06	0.06	0.38	0.52	0.05	-	-	
	CER	my	ENG	0.74	0.99	1.08	2.29	2.50	2.94	0.98	2.74	4.19	2.36	4.94	0.73	1.55	0.98	
			SEA	0.78	0.98	1.19	4.82	9.57	3.33	1.72	1.30	3.24	6.56	8.36	0.75	-	-	
Bud500	WER	tl	ENG	0.14	0.14	0.84	5.60	9.00	3.35	1.19	0.41	0.50	0.23	0.63	0.09	0.12	0.11	
			SEA	0.17	0.17	1.37	3.10	4.06	2.24	3.10	0.41	1.84	0.58	1.48	0.08	-	-	
	WER	vi	ENG	0.10	0.05	0.11	1.50	38.71	1.62	1.45	0.08	0.07	7.09	8.89	0.15	0.47	0.21	
			SEA	0.11	0.34	0.11	3.55	28.94	3.48	2.89	0.08	0.19	6.68	7.18	0.15	-	-	
Commonvoice	WER	en	ENG	0.08	0.09	0.13	0.08	0.07	0.10	0.13	0.08	0.07	0.23	0.38	0.11	0.09	0.10	
			SEA	0.08	0.09	0.13	0.08	0.07	0.10	0.13	0.08	0.07	0.23	0.38	0.11	-	-	
	WER	id	ENG	0.11	0.11	0.10	1.44	0.68	0.38	0.74	0.11	0.10	0.47	0.25	0.03	0.08	0.06	
			SEA	0.07	0.13	0.19	1.61	0.60	0.37	0.79	0.12	0.10	0.68	1.40	0.03	-	-	
	WER	ta	ENG	0.50	0.52	1.49	1.83	4.34	1.33	1.41	1.13	1.41	0.68	1.95	0.29	0.64	0.46	
			SEA	0.33	0.58	1.41	3.36	2.34	1.37	1.97	1.69	1.13	2.33	2.93	0.28	-	-	
	CER	th	ENG	0.16	0.1	0.05	2.35	1.38	0.53	1.31	0.14	0.76	0.51	4.30	0.03	0.07	0.07	
			SEA	0.09	0.21	0.05	8.41	1.14	0.45	1.27	0.16	0.78	11.41	5.79	0.04	-	-	
	WER	vi	ENG	0.46	0.56	0.48	1.25	0.86	0.78	1.34	0.49	0.48	1.10	1.07	0.09	0.45	0.14	
			SEA	0.16	0.73	0.48	1.72	1.02	0.76	0.11	0.48	0.47	1.28	0.97	0.09	-	-	
	CER	zh	ENG	0.12	0.14	0.10	0.13	0.07	0.47	0.24	0.06	0.05	0.80	2.62	0.24	0.18	0.11	
			SEA	0.12	0.16	0.08	0.07	0.05	0.44	0.12	0.06	0.05	2.14	2.36	0.11	-	-	
ESD	WER	en	ENG	0.05	0.05	0.06	0.03	0.05	0.16	0.08	0.04	0.04	0.14	0.27	0.05	0.04	0.06	
			SEA	0.05	0.05	0.06	0.03	0.05	0.16	0.08	0.04	0.04	0.14	0.27	0.05	-	-	
CER	zh	ENG	0.05	0.05	0.11	0.05	0.02	0.33	0.21	0.02	0.02	0.51	1.76	0.12	0.04	0.05		
		SEA	0.05	0.05	0.07	0.02	0.02	0.39	0.05	0.02	0.02	0.74	1.11	0.08	-	-		
FLEURS	WER	tl	ENG	0.16	0.18	1.06	5.92	0.94	1.41	1.98	0.80	0.57	0.07	0.16	0.07	0.12	0.07	
			SEA	0.16	0.22	1.15	4.34	0.60	1.01	1.65	0.89	1.83	0.17	0.18	0.08	-	-	
	WER	id	ENG	0.06	0.10	0.45	4.07	0.47	0.34	0.74	0.18	0.10	0.54	0.09	0.03	0.07	0.04	
			SEA	0.06	0.12	0.12	3.19	0.36	0.28	0.64	0.11	0.31	0.10	0.09	0.03	-	-	
	CER	km	ENG	0.86	1.71	1.00	4.42	8.21	1.14	1.06	3.23	4.55	1.12	4.90	0.17	1.35	0.31	
			SEA	0.98	1.43	2.16	2.25	4.42	0.92	1.40	11.63	4.48	2.41	7.83	0.20	-	-	
	CER	lo	ENG	0.39	0.86	1.01	2.45	4.52	1.29	1.08	1.39	2.63	0.15	1.42	0.16	1.02	0.38	
			SEA	0.51	0.81	1.36	1.53	6.34	1.24	1.20	1.44	1.80	1.75	2.60	0.22	-	-	
	WER	ms	ENG	0.11	0.14	0.76	3.34	1.32	0.44	1.03	0.28	0.27	0.86	0.18	0.05	0.08	0.06	
			SEA	0.11	0.19	0.29	2.74	0.35	0.42	0.87	0.30	0.17	0.15	0.15	0.05	-	-	
	CER	my	ENG	0.79	1.19	1.10	1.85	10.20	3.62	1.07	2.62	4.43	0.15	3.16	0.13	1.16	0.44	
			SEA	0.76	1.03	1.10	4.35	8.87	3.07	1.11	1.99	4.01	1.42	1.97	0.12	-	-	
CER	th	ENG	0.14	0.13	0.10	4.14	2.93	0.57	1.04	0.16	1.20	0.15	4.03	0.07	0.09	0.04		
		SEA	0.11	0.22	0.10	5.40	2.45	0.49	1.14	0.15	0.95	1.42	4.73	0.04	-	-		
WER	vi	ENG	0.10	0.16	0.91	4.46	2.39	0.40	1.05	0.10	0.09	0.17	0.30	0.06	0.08	0.04		
		SEA	0.08	0.26	0.16	4.11	1.53	0.34	1.42	0.11	0.09	0.19	0.32	0.04	-	-		
CER	zh	ENG	0.10	0.12	0.84	0.1	0.08	0.36	0.18	0.07	0.06	0.25	1.00	0.18	0.08	0.07		
		SEA	0.10	0.11	0.84	0.07	0.08	0.35	0.10	0.07	0.06	1.00	1.00	0.09	-	-		
ASR-MALCSC	WER	ms	ENG	0.21	0.26	0.79	3.82	14.23	1.30	1.77	0.21	1.01	4.86	2.81	0.26	0.31	0.42	
			SEA	0.22	0.27	0.71	2.66	7.70	1.84	1.68	0.71	1.24	4.07	4.58	0.26	-	-	
MIG	CER	my	ENG	1.06	1.59	1.47	2.58	12.47	8.30	1.47	1.12	9.45	8.13	1.63	0.31	2.55	0.78	
			SEA	0.97	1.53	2.09	9.35	15.58	2.80	2.09	1.64	4.43	3.58	5.53	0.29	-	-	
OpenSLR	CER	km	ENG	0.68	1.72	1.14	1.80	2.33	2.29	1.10	2.71	1.61	2.29	2.94	0.09	4.70	0.29	
			SEA	0.88	1.74	2.18	5.91	4.64	2.80	1.13	8.42	1.70	2.89	3.92	0.30	-	-	
	CER	my	ENG	0.71	0.97	1.16	1.80	6.78	5.34	1.16	1.57	3.88	1.26	1.51	0.05	2.51	0.42	
			SEA	0.74	1.07	2.63	7.51	13.1	5.52	2.13	0.97	4.40	0.99	1.42	0.06	-	-	
WER	ta	ENG	0.26	0.34	1.55	2.03	4.58	1.06	1.55	1.58	1.26	0.22	0.45	0.24	2.51	0.35		
		SEA	0.26	0.40	1.30	3.59	2.47	1.10	2.22	1.58	1.47	0.71	0.85	0.25	-	-		
SFDUSC	WER	tl	ENG	0.24	0.27	1.20	3.77	0.31	2.44	2.48	0.59	0.59	0.42	0.48	0.23	0.25	0.24	
			SEA	0.25	0.10	1.11	3.82	1.30	2.24	4.07	0.70	0.61	0.39	0.75	0.21	-	-	
ASR-SgpCCSC	CER	zh	ENG	0.07	0.11	0.30	0.13	0.04	0.62	0.19	0.05	0.05	1.19	0.80	0.19	1.00	0.13	
			SEA	0.07	0.09	0.21	3.82	0.04	0.71	0.21	0.05	0.04	5.55	1.77	0.12	-	-	
SG Streets	WER	en	ENG	0.10	0.05	0.57	0.17	0.82	0.39	0.20	0.09	0.09	2.27	4.11	0.22	1.00	0.29	
			SEA	0.10	0.05	0.57	0.17	0.82	0.39	0.20	0.09	0.09	2.27	4.11	0.22	-	-	
ASR-SMalDuSC	WER	ms	ENG	0.07	0.11	0.27	2.28	2.21	0.55	0.27	0.21	0.22	0.23	0.18	0.02	1.00	0.03	
			SEA	0.07	0.10	0.27	2.24	1.30	0.56	1.08	0.21	0.23	0.26	0.24	0.02	-	-	
Thai Elderly Speech	CER	th	ENG	0.08	0.07	0.04	2.07	1.56	0.71	1.35	0.14	0.17	0.54	2.72	0.07	0.06	0.10	
			SEA	0.08	0.22	0.04	8.38	1.16	0.75	1.32	0.20	0.12	7.46	2.63	0.05	-	-	
	LOTUS	CER	th	ENG	0.02	0.03	0.02	3.35	4.51	0.46	1.14	0.04	0.03	0.06	1.06	0.01	0.03	0.01
				SEA	0.02	0.12	0.03	6.51	1.37	0.50	0.03	0.04	0.03	1.25	1.17	0.03	-	-
VietMed	WER	vi	ENG	0.27	0.64	0.27	4.47	21.47	0.92	0.27	0.40	0.35	2.49	4.80	0.18	0.65	0.44	
			SEA	0.28	0.49	0.27	3.23	19.44	1.48	1.56	0.40	0.26	2.71	4.07	0.19	-	-	
Average		ENG		0.27	0.41	0.65	2.29	4.85	1.42	0.95	0.67	1.24	1.27	2.01	0.14	0.74	0.22	
			SEA	0.27	0.43	0.74	3.39	4.30	1.45	1.20	1.11	1.08	2.28	2.50	0.14	-	-	

Table 12: Age Recognition performance: results across datasets comparing English and SEA prompts. Evaluation is conducted using the Model-as-Judge metric, where higher scores indicate better performance.

Model			MERaLiON2 10B	MERaLiON2 3B	SeaLLMs Audio 7B	Phi-4 multi- -modal instruct	Kimi Audio	Voxtral mini	Qwen2 Audio 7B	Qwen 2.5 Omni 3B	Qwen 2.5 Omni 7B	gemma 3n E4B-it	gemma 3n E2B-it	Gemini 2.5 Flash	GPT 4o Audio
Size			10B	3B	7B	5.6B	7B	3B	7B	3B	7B	4B	2B	-	-
Data	Lang Prompt														
Commonvoice	en	ENG	63.10	61.00	54.10	45.35	38.00	58.10	36.90	29.30	36.90	63.60	59.45	72.00	65.00
		SEA	63.10	61.00	54.10	45.35	38.00	58.10	36.90	29.30	36.90	63.60	59.45	72.00	65.00
	ta	ENG	64.65	71.60	52.70	49.20	51.70	77.00	18.00	48.50	18.00	65.50	63.95	75.00	71.00
		SEA	47.90	16.10	1.10	16.30	14.80	38.10	19.35	7.00	19.35	38.40	17.00	77.00	62.00
	th	ENG	57.81	67.23	68.77	40.39	49.10	81.81	16.90	48.52	16.90	73.42	72.13	74.00	78.00
		SEA	42.19	51.87	1.00	53.16	63.48	61.42	0.25	43.74	0.26	68.13	49.16	85.00	53.00
	vi	ENG	73.23	66.57	63.03	34.81	55.52	85.83	6.96	58.58	6.96	71.79	66.03	86.00	86.00
		SEA	64.35	52.34	1.00	30.09	62.61	33.73	1.56	33.13	1.56	61.40	33.01	90.00	81.00
zh	ENG	72.40	74.60	77.20	47.50	45.25	75.80	22.30	67.60	22.00	74.10	68.00	73.00	83.00	
	SEA	69.00	51.00	1.00	65.00	63.20	71.70	57.60	59.20	57.60	62.15	35.20	73.00	45.00	
Average	ENG		66.24	68.20	63.16	43.45	47.91	75.71	20.21	50.50	20.15	69.68	65.91	76.00	76.60
	SEA		57.31	46.46	11.64	41.98	48.42	52.61	23.13	34.47	23.13	58.74	38.76	79.40	61.20

Table 13: Emotion Recognition performance: results across datasets comparing English and SEA prompts. Evaluation is conducted using the Model-as-Judge metric, where higher scores indicate better performance.

Model			MERaLiON2 10B	MERaLiON2 3B	SeaLLMs Audio 7B	Phi-4 multi- -modal instruct	Kimi Audio	Voxtral mini	Qwen2 Audio 7B	Qwen 2.5 Omni 3B	Qwen 2.5 Omni 7B	gemma 3n E4B-it	gemma 3n E2B-it	Gemini 2.5 Flash	GPT 4o Audio
Size			10B	3B	7B	5.6B	7B	3B	7B	3B	7B	4B	2B	-	-
Data	Lang Prompt														
EmoTa	ta	ENG	8.33	15.12	6.78	21.58	15.33	8.33	23.08	7.26	14.58	8.49	8.23	12.00	8.00
		SEA	11.97	10.84	0.85	2.24	7.91	3.73	0.96	0.85	1.00	6.94	11.43	9.00	17.00
ESD	en	ENG	19.50	23.40	12.85	19.80	64.70	9.05	40.90	9.40	13.55	8.50	9.15	15.00	13.00
		SEA	19.50	23.40	12.85	19.80	64.70	9.05	40.90	9.40	13.55	8.50	9.15	15.00	13.00
		zh	ENG	16.70	19.90	7.70	17.80	71.25	5.55	34.60	15.50	18.10	10.55	11.60	14.00
IndoWaveSentiment	id	SEA	14.15	18.75	7.45	12.50	66.15	4.80	47.55	17.70	16.10	12.65	13.65	14.00	13.50
		ENG	20.00	20.67	11.67	18.33	27.83	11.00	13.33	11.33	13.33	17.00	15.83	26.00	23.00
M3ED	zh	SEA	23.33	18.00	17.50	3.67	29.83	8.67	15.50	16.00	14.33	13.33	9.67	18.00	22.00
		ENG	19.35	26.65	10.90	20.65	19.00	7.95	14.30	11.95	13.50	7.35	10.20	16.00	12.00
TEC	ta	SEA	23.25	23.40	13.35	13.50	23.15	0.85	19.80	18.55	19.60	13.20	10.75	16.50	14.00
		ENG	34.85	42.42	24.85	30.30	42.42	23.12	30.91	28.48	27.58	21.52	19.39	42.50	43.00
THAI SER	th	SEA	36.97	20.30	1.82	3.64	21.81	6.06	2.42	1.21	1.21	31.52	27.88	35.00	47.00
		ENG	12.36	19.74	11.62	17.59	17.48	9.31	14.19	10.26	13.66	13.82	11.10	11.00	13.00
Average	ENG		22.56	26.55	15.53	22.99	35.88	14.55	26.74	15.79	19.10	17.61	16.75	23.81	19.88
	SEA		23.97	21.96	12.76	12.78	40.86	9.94	22.26	12.72	13.95	18.90	17.59	21.44	22.06

Table 14: Gender Recognition performance: results across datasets comparing English and SEA prompts. Evaluation is conducted using the Model-as-Judge metric, where higher scores indicate better performance.

Model		MERaLiON2 10B	MERaLiON2 3B	SeaLLMs Audio 7B	Phi-4 multi- modal instruct	Kimi Audio	Voxtral mini	Qwen2 Audio 7B Instruct	Qwen 2.5 Omni 3B	Qwen 2.5 Omni 7B	gemma 3n E4B-it	gemma 3n E2B-it	Gemini 2.5 Flash	GPT 4o Audio	
Size		10B	3B	7B	5.6B	7B	3B	7B	3B	7B	4B	2B	-	-	
Data	Lang Prompt														
Commonvoice	id	ENG	45.20	45.30	54.90	35.80	93.50	32.00	95.20	41.20	59.30	49.60	22.20	95.00	46.00
		SEA	57.30	62.92	56.18	57.30	66.30	44.94	61.80	62.92	52.81	57.30	60.67	82.02	53.93
	ta	ENG	53.00	51.40	50.10	45.70	90.60	44.60	96.80	23.80	53.90	54.70	37.20	95.00	33.00
		SEA	40.40	46.90	0.10	7.90	57.60	13.20	7.40	0.40	4.10	12.30	17.20	92.00	35.00
	th	ENG	50.07	36.14	57.43	28.51	96.12	44.18	96.79	46.18	61.45	51.54	29.18	93.00	50.00
		SEA	23.96	53.55	49.00	28.92	89.42	3.88	96.39	32.33	67.47	19.14	2.54	91.00	40.00
	vi	ENG	24.05	18.82	83.01	34.25	94.51	50.07	95.16	42.75	49.41	38.04	33.73	90.00	26.00
		SEA	14.64	43.14	81.57	18.30	94.41	22.09	97.52	21.44	9.80	3.66	3.53	77.00	35.00
zh	ENG	53.70	34.60	68.90	57.10	83.90	40.90	98.20	75.30	81.80	41.00	17.80	90.00	49.00	
	SEA	35.50	36.50	68.90	44.70	64.50	31.40	98.40	60.90	88.90	37.95	26.60	91.00	21.00	
EmoTa	ta	ENG	67.31	49.36	52.88	40.81	90.92	22.01	98.82	17.41	36.81	55.34	49.36	94.00	25.00
		SEA	48.93	47.76	0.21	8.65	75.43	10.47	2.35	1.28	2.67	11.65	9.19	94.00	33.00
FLEURS	en	ENG	58.27	67.54	39.41	62.44	97.84	8.04	99.38	44.20	30.76	52.24	26.35	97.00	78.00
		SEA	58.27	67.54	39.41	62.44	97.84	8.04	99.38	44.20	30.76	52.24	26.35	97.00	78.00
	km	ENG	56.60	31.90	72.03	29.93	99.22	38.04	78.82	33.73	43.66	48.89	11.37	99.00	62.00
		SEA	43.40	56.34	0.39	26.08	47.06	1.96	5.49	2.35	1.05	2.75	0.78	97.00	15.00
IndoWaveSentiment	id	ENG	71.67	63.67	50.33	42.00	91.33	5.33	98.33	13.00	65.00	58.00	20.00	96.00	60.00
		SEA	60.67	3.60	50.00	22.33	79.33	1.33	94.00	19.33	31.67	24.00	13.33	95.00	14.00
M3ED	zh	ENG	84.30	64.50	51.40	82.80	74.60	8.80	86.30	15.80	43.30	58.05	30.80	84.00	23.00
		SEA	70.70	80.60	53.80	72.10	78.90	6.10	78.00	11.30	78.90	45.45	24.20	81.00	12.00
OpenSLR	ta	ENG	55.30	50.50	50.30	47.00	97.50	34.00	99.10	41.10	40.50	52.00	29.10	97.00	47.00
		SEA	37.80	42.20	0.10	6.20	60.80	9.40	8.00	0.80	4.70	15.10	14.60	97.00	36.00
SG Streets	en	ENG	89.63	39.43	61.99	53.46	95.94	5.49	84.55	8.94	5.69	64.23	31.30	98.00	32.00
		SEA	89.63	39.43	61.99	53.46	95.94	5.49	84.55	8.94	5.69	64.23	31.30	98.00	32.00
ASR-SMaIDuSC	ms	ENG	52.40	49.80	52.20	22.30	96.70	26.80	93.90	66.50	58.90	49.70	28.00	98.00	76.00
		SEA	44.00	50.30	45.70	37.00	94.30	19.00	61.30	23.55	15.90	34.70	20.20	96.00	24.00
Thai Elderly Speech	th	ENG	68.15	55.75	34.48	22.58	94.96	44.76	97.48	45.16	60.28	65.73	26.66	99.00	46.00
		SEA	26.92	46.47	69.56	45.67	84.68	3.33	96.98	23.99	82.01	37.40	2.42	85.00	51.00
THAI SER	th	ENG	63.46	63.87	49.21	53.40	86.81	18.32	87.54	31.10	60.00	56.86	27.23	86.00	44.00
		SEA	61.78	58.12	81.15	41.68	84.29	12.14	85.86	26.13	80.00	33.93	3.77	85.00	34.00
Vietnam-Celeb	vi	ENG	65.80	54.80	50.90	49.30	71.10	19.10	69.70	31.70	56.60	57.80	23.10	69.00	41.00
		SEA	61.40	63.50	50.30	30.30	72.80	5.00	69.30	23.00	27.50	10.65	1.60	73.00	36.00
Average	ENG	SEA	59.78	49.43	55.04	44.98	89.52	28.67	90.46	37.69	50.60	53.59	29.65	91.88	46.58
		SEA	46.93	47.33	44.40	34.05	78.47	11.95	66.05	22.70	35.65	29.00	17.36	89.59	32.94

Table 15: Speaker Recognition performance: results across datasets comparing English and SEA prompts. Evaluation is conducted using the Model-as-Judge metric, where higher scores indicate better performance.

Model			MERaLiON2 10B	MERaLiON2 3B	SeaLLMs Audio 7B	Phi-4 multi- -modal instruct	Kimi Audio	Voxtral mini	Qwen2 Audio 7B	Qwen 2.5 Omni 3B	Qwen 2.5 Omni 7B	gemma 3n E4B-it	gemma 3n E2B-it	Gemini 2.5 Flash	GPT 4o Audio
Size	10B	3B	7B	5.6B	7B	3B	7B	3B	7B	3B	7B	4B	2B	-	-
Data	Lang	Prompt													
EmoTa	ta	ENG	53.10	43.06	56.94	14.10	57.05	46.15	50.64	26.50	13.25	39.53	41.99	71.00	6.00
		SEA	53.85	31.73	0.00	15.71	51.17	33.97	15.44	7.91	12.39	27.99	40.92	73.00	13.00
ESD	en	ENG	52.90	39.20	56.75	34.90	53.80	37.70	47.90	24.40	5.70	39.10	45.90	71.00	10.00
		SEA	52.90	39.20	56.75	34.90	53.80	37.70	47.90	24.40	5.70	39.10	45.90	71.00	10.00
	zh	ENG	53.10	43.70	49.10	24.80	48.00	42.40	51.60	31.70	23.30	41.40	41.20	59.00	12.00
		SEA	50.60	43.30	48.10	45.20	44.90	43.50	45.30	38.10	30.70	45.58	44.40	61.00	4.00
mig	my	ENG	50.00	35.15	53.50	22.50	50.90	31.50	54.60	30.60	13.50	35.90	29.40	70.00	5.00
		SEA	36.50	12.70	0.20	11.60	12.60	27.20	1.00	6.30	6.45	43.80	37.30	77.00	4.00
ASR-SMaDuSC	ms	ENG	59.20	44.50	57.80	23.20	59.30	47.00	46.60	24.00	10.80	37.40	44.00	74.00	14.00
		SEA	48.05	39.30	42.20	28.20	59.70	44.70	41.10	29.40	17.20	47.60	38.90	82.00	16.00
Thai Elderly Speech	th	ENG	52.92	38.13	52.29	20.42	52.10	50.31	48.02	29.48	14.27	42.81	43.65	75.00	6.00
		SEA	34.06	26.77	37.92	27.40	56.25	41.56	43.65	6.56	5.00	38.49	30.52	71.00	7.00
THAI SER	th	ENG	51.02	38.45	50.70	22.13	58.00	46.72	50.16	30.93	18.80	41.57	46.94	70.00	20.00
		SEA	42.96	36.84	42.86	30.08	60.47	34.80	42.11	12.24	6.20	37.16	35.02	61.00	7.00
VoxVietnam-O	vi	ENG	53.80	35.80	51.60	24.10	54.10	41.40	49.30	31.10	17.40	41.30	45.70	79.00	11.00
		SEA	50.00	39.80	46.90	26.80	46.40	25.30	56.00	17.40	7.90	38.00	41.60	78.00	6.00
Average	ENG		53.25	39.75	53.59	23.27	54.16	42.90	49.85	28.59	14.63	39.88	42.35	71.13	10.50
	SEA		46.12	33.71	34.37	27.48	48.16	36.09	36.56	17.79	11.44	39.72	39.32	71.75	8.38

Table 16: Speech Translation performance: results across datasets comparing English and SEA prompts. Evaluation is based on the BLEU metric, where higher scores indicate better performance.

Model			MERaLiON2 10B	MERaLiON2 3B	SeaLLMs Audio 7B	Phi-4 multi- -modal instruct	Kimi Audio	Voxtral mini	Qwen2 Audio 7B Instruct	Qwen 2.5 Omni 3B	Qwen 2.5 Omni 7B	gemma 3n E4B-it	gemma 3n E2B-it	Gemini 2.5 Flash	GPT 4o Audio
Size			10B	3B	7B	5.6B	7B	3B	7B	3B	7B	4B	2B	-	-
Data	Lang	Prompt													
FLEURS	id	ENG	32.48	20.97	25.55	0.70	10.75	35.02	9.37	16.65	16.55	21.40	20.52	24.07	33.44
		SEA	30.02	21.98	25.97	0.13	16.81	34.58	6.23	11.15	17.93	24.80	17.57	24.43	34.90
	km	ENG	2.19	0.50	0.73	0.23	0.35	7.04	0.47	0.26	0.48	4.72	2.52	12.40	6.34
		SEA	1.97	0.45	0.33	0.05	0.49	0.08	0.11	0.27	0.24	4.71	2.50	13.76	7.52
	lo	ENG	11.34	1.03	5.91	0.04	0.55	10.44	0.47	2.49	3.06	7.66	6.75	15.36	13.61
		SEA	12.09	0.88	1.04	0.01	1.04	0.55	0.11	0.80	2.35	8.94	7.18	15.04	13.06
	ms	ENG	31.04	15.52	19.71	1.04	6.15	30.44	5.25	12.16	13.19	18.24	16.41	26.49	33.53
		SEA	34.35	13.09	20.30	0.22	20.02	30.33	2.57	10.12	14.24	23.04	14.37	23.70	33.93
	my	ENG	0.37	0.12	0.29	0.11	0.03	0.73	0.40	0.11	0.35	0.78	0.35	8.59	1.47
		SEA	0.59	0.14	0.03	0.11	0.04	0.63	0.14	0.02	0.10	0.73	0.10	15.05	2.08
	th	ENG	17.22	3.47	12.30	0.14	2.42	19.53	0.47	7.90	8.34	12.95	9.57	16.93	22.41
		SEA	18.70	5.59	11.51	0.02	3.52	20.15	0.55	7.86	9.29	14.90	8.49	20.25	22.93
	tl	ENG	25.79	8.28	1.67	0.92	5.99	30.42	1.68	2.52	2.50	17.14	12.98	14.50	28.48
		SEA	27.22	7.76	1.76	0.10	23.27	30.69	1.25	0.89	2.29	22.47	14.44	16.09	26.20
	vi	ENG	19.18	6.50	13.52	0.12	3.94	24.79	2.05	11.40	12.02	7.19	4.89	17.26	27.10
		SEA	22.91	5.95	13.00	0.01	4.19	25.05	1.91	10.31	12.32	11.84	7.51	19.74	30.27
	zh	ENG	20.14	11.63	17.01	24.08	0.01	21.36	20.69	14.89	14.73	8.72	6.75	16.18	24.77
		SEA	20.74	10.94	17.11	2.28	0.01	21.30	15.88	15.02	13.81	10.78	6.37	21.91	21.60
Average	ENG		17.75	7.56	10.74	3.04	3.36	19.98	4.54	7.60	7.91	10.98	8.97	16.86	21.24
	SEA		18.73	7.42	10.11	0.32	7.71	18.15	3.20	6.27	8.06	13.58	8.72	18.89	21.39

Table 17: TCQ performance: results across datasets comparing English and SEA prompts. Scores are reported as raw WER and CER values without normalization; lower values indicate better performance.

Model				MERaLiON2	MERaLiON2	SeaLLMs	Phi-4	Kimi	Voxtral	Qwen2	Qwen	Qwen	gemma	gemma	Gemini	GPT		
Size				10B	3B	Audio	multi-	Audio	mini	Audio	2.5	2.5	3n	3n	2.5	4o		
Data Lang Length (Max) Prompt				10B	3B	7B	modal	7B	3B	7B	Omni	Omni	E4B-it	E2B-it	Flash	Audio		
							instruct			Instruct	3B	7B	4B	2B	-	-		
SG Streets	en	30	ENG	2.43	2.43	2.12	2.67	2.47	1.89	2.30	2.64	4.21	2.20	2.32	0.47	2.24		
			SEA	2.43	2.43	2.12	2.67	2.47	1.89	2.30	2.64	4.21	2.20	2.32	0.47	2.24		
	60	ENG	3.44	3.47	-	2.10	3.36	2.12	-	2.29	3.35	1.91	1.94	0.76	2.49			
		SEA	3.44	3.47	-	2.10	3.36	2.12	-	2.29	3.35	1.91	1.94	0.76	2.49			
		ASR-SgpCCSC	zh	30	ENG	4.78	5.44	6.54	5.48	9.94	6.84	9.08	3.91	3.69	8.84	9.55	1.57	5.17
		SEA	4.34	5.13	5.68	4.07	10.43	5.46	8.87	4.15	2.91	8.62	9.67	1.57	4.99			
YODAS2	en	30	ENG	9.85	11.28	-	6.72	22.74	12.86	-	3.85	5.73	9.98	9.36	1.89	8.00		
			SEA	9.50	10.25	-	4.03	20.95	8.16	-	3.77	3.71	9.08	8.83	1.89	9.21		
	60	ENG	20.93	21.15	-	18.92	41.07	27.07	-	-	18.59	-	-	3.10	8.68			
		SEA	20.86	17.09	-	14.94	37.36	12.12	-	-	8.19	-	-	3.10	10.81			
		ENG	20.59	-	-	-	57.41	36.85	-	-	-	-	-	3.30	7.35			
		SEA	21.81	-	-	-	56.92	16.92	-	-	-	-	-	3.30	9.04			
id	30	ENG	5.34	5.31	3.81	4.53	5.79	3.62	4.57	4.94	4.42	5.03	5.03	1.36	4.42			
		SEA	5.34	5.31	3.81	4.53	5.79	3.62	4.57	4.94	4.42	5.03	5.03	1.36	4.42			
	60	ENG	10.98	10.74	-	6.39	11.76	5.48	-	7.01	4.86	5.83	5.96	1.29	6.37			
		SEA	10.98	10.74	-	6.39	11.76	5.48	-	7.01	4.86	5.83	5.96	1.29	6.37			
	120	ENG	16.19	15.10	-	9.02	22.28	9.09	-	-	6.92	-	-	1.82	6.49			
		SEA	16.19	15.10	-	9.02	22.28	9.09	-	-	6.92	-	-	1.82	6.49			
	180	ENG	-	-	-	-	-	-	-	-	-	-	-	1.75	6.72			
		SEA	-	-	-	-	-	-	-	-	-	-	-	1.75	6.72			
	th	30	ENG	3.82	3.70	3.42	21.41	22.32	2.71	4.05	3.76	3.23	3.78	3.81	0.64	3.18		
			SEA	3.77	3.50	3.52	14.13	9.58	3.06	3.86	3.81	4.54	4.00	3.91	1.36	3.20		
		60	ENG	7.88	6.44	-	14.18	19.97	4.40	-	5.26	4.09	4.23	4.30	1.21	4.47		
			SEA	7.63	6.00	-	9.43	8.76	4.75	-	5.53	6.53	4.24	4.44	1.79	5.23		
120		ENG	13.07	9.45	-	13.92	16.49	7.27	-	-	6.09	-	-	1.38	4.88			
		SEA	12.59	9.08	-	12.35	10.94	6.27	-	-	8.20	-	-	2.37	4.89			
vi	30	ENG	14.00	-	-	-	16.24	10.26	-	-	-	-	-	2.02	5.42			
		SEA	13.86	-	-	-	9.26	7.29	-	-	-	-	-	2.41	4.94			
	60	ENG	3.82	1.49	5.31	80.08	27.81	5.24	3.14	36.29	9.68	9.21	10.49	10.15	7.09			
		SEA	3.75	1.57	6.33	40.13	16.43	6.93	2.47	17.47	12.56	9.85	13.22	8.17	8.11			
	120	ENG	6.05	1.77	-	51.66	37.37	8.21	-	36.06	11.39	9.79	9.80	9.32	10.31			
		SEA	5.63	1.73	-	37.39	17.88	9.97	-	23.94	13.33	9.74	9.89	7.43	11.61			
180	ENG	10.10	2.10	-	41.25	32.92	12.71	-	-	11.80	-	-	10.16	11.10				
	SEA	8.87	1.86	-	37.47	19.98	14.27	-	-	14.77	-	-	8.02	11.84				
zh	30	ENG	10.52	-	-	-	60.78	14.44	-	-	-	-	-	12.21	12.98			
		SEA	9.76	-	-	-	27.66	16.71	-	-	-	-	-	9.93	12.87			
	60	ENG	5.48	5.22	4.11	21.25	17.42	4.28	4.42	3.06	3.11	5.30	5.95	1.26	4.33			
		SEA	5.46	4.92	4.32	14.75	19.24	4.72	4.68	3.50	5.37	5.46	6.21	1.39	4.29			
	120	ENG	10.60	7.46	-	12.86	15.43	7.54	-	3.80	3.78	5.79	5.85	1.12	5.16			
		SEA	10.42	6.25	-	11.12	17.28	7.88	-	3.73	6.13	5.89	5.84	1.46	5.90			
180	ENG	14.05	8.90	-	15.76	11.24	11.77	-	-	4.81	-	-	46.66	5.90				
	SEA	13.95	7.34	-	19.65	23.72	11.90	-	-	6.45	-	-	2.01	5.72				
en	30	ENG	12.22	-	-	-	17.52	16.12	-	-	-	-	-	1.77	5.31			
		SEA	12.03	-	-	-	26.28	12.47	-	-	-	-	-	2.06	4.55			
	60	ENG	10.16	9.83	13.04	9.29	12.21	8.61	11.38	6.57	10.10	12.95	14.12	1.41	8.99			
		SEA	8.46	9.40	11.19	8.49	12.23	6.91	11.00	5.31	7.93	12.94	13.58	2.98	10.74			
	120	ENG	17.60	16.31	-	10.41	22.87	13.54	-	6.15	12.89	14.16	15.26	2.86	10.92			
		SEA	15.95	14.34	-	11.48	23.07	8.33	-	4.67	9.09	13.88	15.08	4.40	13.01			
180	ENG	28.66	22.18	-	14.92	44.31	22.13	-	-	20.73	-	-	2.94	14.59				
	SEA	25.32	19.94	-	11.47	43.27	10.74	-	-	6.88	-	-	3.73	15.04				
en	30	ENG	31.04	-	-	-	61.92	33.00	-	-	-	-	-	3.51	13.76			
		SEA	28.22	-	-	-	73.53	14.62	-	-	-	-	-	4.49	17.39			
	60	ENG	17.67	-	-	-	42.78	22.13	-	-	-	-	-	4.09	8.59			
		SEA	17.14	-	-	-	38.73	13.60	-	-	-	-	-	3.99	9.25			
	120	ENG	14.72	11.27	-	16.25	24.04	12.86	-	-	9.85	-	-	9.44	7.38			
		SEA	13.97	10.06	-	14.99	22.51	9.20	-	-	7.34	-	-	3.01	7.83			
30	ENG	5.12	4.77	5.48	20.67	14.00	4.74	5.56	8.74	5.49	6.76	7.32	2.41	5.06				
	SEA	4.79	4.61	5.28	12.68	10.88	4.66	5.39	5.97	5.99	6.87	7.70	2.47	5.43				
60	ENG	9.48	8.21	-	14.90	19.07	7.74	-	9.20	6.58	7.38	7.49	2.64	6.82				
	SEA	9.08	7.54	-	11.70	14.72	6.67	-	7.28	6.71	7.22	7.42	2.72	7.69				

Table 18: TLoc performance: results across datasets comparing English and SEA prompts. Evaluation is based on the F1 score metric, where higher values indicate better performance.

Model				MERaLiON2	MERaLiON2	SeaLLMs	Phi-4	Kimi	Voxtral	Qwen2	Qwen	Qwen	gemma	gemma	Gemini	GPT
Size				10B	3B	7B	multi-	Audio	mini	Audio	2.5	2.5	3n	3n	2.5	4o
Prompt				10B	3B	7B	modal			Instruct	Omni	Omni	E4B-it	E2B-it	Flash	Audio
Data	Lang	Length (Max)	Prompt				instruct	7B	3B	7B	3B	7B	4B	2B	-	-
ASR-SgpCCSC	zh	30	ENG	11.10	10.85	8.45	19.90	-	-	19.10	13.10	31.21	7.86	6.44	11.89	22.19
			SEA	11.10	10.85	8.45	19.90	-	-	19.10	13.10	31.21	7.86	6.44	11.89	22.19
	60	ENG	6.23	6.42	-	7.73	-	-	-	10.78	12.71	6.64	2.88	1.26	10.50	
		SEA	6.23	6.42	-	7.73	-	-	-	10.78	12.71	6.64	2.88	1.26	10.50	
		ENG	4.18	2.59	-	4.66	-	-	-	-	7.84	-	-	4.55	6.14	
		SEA	4.18	2.59	-	4.66	-	-	-	-	7.84	-	-	4.55	6.14	
SG Streets	en	30	ENG	40.28	33.56	22.56	11.96	17.17	27.70	54.85	49.61	44.77	18.04	19.18	16.62	40.25
			SEA	40.28	33.56	22.56	11.96	17.17	27.70	54.85	49.61	44.77	18.04	19.18	16.62	40.25
	60	ENG	23.75	15.16	-	5.30	11.13	17.90	-	29.78	38.05	7.36	12.87	17.90	37.94	
	SEA	23.75	15.16	-	5.30	11.13	17.90	-	29.78	38.05	7.36	12.87	17.90	37.94		
YODAS2	en	30	ENG	22.99	20.00	9.92	12.70	14.78	15.46	41.62	34.01	36.27	16.49	16.35	5.88	24.30
			SEA	22.99	20.00	9.92	12.70	14.78	15.46	41.62	34.01	36.27	16.49	16.35	5.88	24.30
		60	ENG	13.91	11.46	-	6.79	8.43	7.76	-	21.07	18.67	10.44	11.36	7.98	17.14
			SEA	13.91	11.46	-	6.79	8.43	7.76	-	21.07	18.67	10.44	11.36	7.98	17.14
			ENG	8.28	6.12	-	3.47	2.72	2.62	-	-	10.78	-	-	2.99	7.65
			SEA	8.28	6.12	-	3.47	2.72	2.62	-	-	10.78	-	-	2.99	7.65
	id	30	ENG	27.44	24.58	16.41	14.02	17.06	20.90	38.26	40.91	45.27	15.55	12.60	11.83	28.48
			SEA	26.70	18.43	19.78	10.63	15.04	18.26	37.56	21.84	34.71	20.16	16.19	16.19	30.01
		60	ENG	14.90	12.99	-	7.47	9.94	11.76	-	23.78	21.44	10.50	9.66	11.04	20.84
		SEA	14.82	10.85	-	7.66	8.15	10.77	-	12.86	0.00	14.53	11.59	15.57	22.20	
		ENG	10.13	8.58	-	4.96	5.05	5.47	-	-	16.59	-	-	10.01	15.04	
		SEA	10.43	6.45	-	5.18	4.91	4.17	-	-	13.73	-	-	9.42	15.60	
	th	30	ENG	17.79	17.42	8.94	10.48	12.97	14.24	16.47	25.52	37.62	12.45	9.48	12.83	26.66
			SEA	16.04	15.32	9.16	13.75	13.02	10.84	11.92	14.08	27.42	10.71	6.93	25.73	17.62
		60	ENG	9.03	8.56	-	5.90	8.03	7.16	-	14.76	15.22	9.28	7.72	4.35	18.64
		SEA	9.92	5.08	-	7.48	5.87	6.73	-	9.77	17.20	7.36	5.62	16.38	15.51	
		ENG	5.65	4.13	-	3.04	4.08	3.78	-	-	11.17	-	-	8.16	11.31	
		SEA	5.87	1.94	-	3.70	5.80	3.00	-	-	6.75	-	-	15.81	9.18	
vi	30	ENG	22.40	14.50	8.22	10.30	15.84	18.81	35.08	33.18	39.17	12.14	10.35	10.11	25.67	
		SEA	24.76	13.49	14.04	14.63	8.16	15.44	28.60	31.59	40.73	12.23	10.82	8.70	28.80	
	60	ENG	12.18	11.54	-	4.29	9.76	8.86	-	17.98	24.73	9.35	9.14	6.35	11.98	
	SEA	14.77	7.92	-	7.70	11.35	7.30	-	17.96	28.44	10.95	8.88	7.28	21.40		
	ENG	6.34	7.34	-	2.44	4.84	3.76	-	-	12.73	-	-	2.83	7.40		
	SEA	6.88	6.76	-	3.11	1.20	4.42	-	-	20.39	-	-	3.40	12.68		
zh	30	ENG	13.56	10.80	6.49	11.41	9.13	9.97	27.69	17.10	15.84	10.19	8.32	12.43	17.77	
		SEA	12.52	10.61	6.69	11.02	11.41	11.14	26.21	13.93	16.46	12.36	7.29	10.11	14.80	
	60	ENG	6.61	5.85	-	6.28	4.36	5.76	-	9.35	9.07	7.45	4.29	5.02	7.73	
	SEA	6.80	7.91	-	5.62	5.80	5.61	-	8.50	11.54	6.95	7.15	4.13	4.64		
	ENG	3.85	2.08	-	3.26	1.81	2.93	-	-	8.81	-	-	4.89	3.64		
	SEA	3.98	4.38	-	3.16	0.00	2.35	-	-	4.77	-	-	6.32	4.34		