

WenetSpeech-Wu: Datasets, Benchmarks, and Models for a Unified Chinese Wu Dialect Speech Processing Ecosystem

Anonymous ACL submission

Abstract

Speech processing for low-resource dialects remains a fundamental challenge in developing inclusive and robust speech technologies. Despite its linguistic significance and large speaker population, the Wu dialect of Chinese has long been hindered by the lack of large-scale speech data, standardized evaluation benchmarks, and publicly available models. In this work, we present **WenetSpeech-Wu**, the first large-scale, multi-dimensionally annotated open-source speech corpus for the Wu dialect, comprising approximately 8,000 hours of diverse speech data. Building upon this dataset, we introduce **WenetSpeech-Wu-Bench**, the first standardized and publicly accessible benchmark for systematic evaluation of Wu dialect speech processing, covering automatic speech recognition (ASR), Wu-to-Mandarin translation, speaker attribute prediction, speech emotion recognition, text-to-speech (TTS) synthesis, and instruction-following TTS (instruct TTS). Furthermore, we release a suite of strong open-source models trained on WenetSpeech-Wu, establishing competitive performance across multiple tasks and empirically validating the effectiveness of the proposed dataset. Together, these contributions lay the foundation for a comprehensive Wu dialect speech processing ecosystem, and we open-source proposed datasets, benchmarks, and models to support future research on dialectal speech intelligence ¹.

1 Introduction

Speech processing has become a fundamental component of artificial intelligence, enabling natural and efficient human-machine interaction across a wide range of real-world applications (Peng et al., 2025; Xu et al., 2025a). For high-resource languages such as Mandarin Chinese and English, the speech processing ecosystem has reached a high

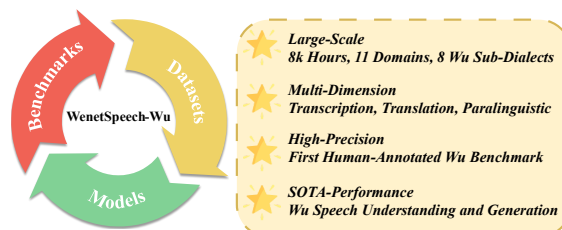


Figure 1: Highlights of WenetSpeech-Wu.

level of maturity, characterized by large-scale and diverse datasets (Zhang et al., 2022; Zen et al., 2019; He et al., 2024), publicly available and multi-dimensional evaluation benchmarks (Sakshi et al., 2025; Wang et al., 2025a,b), and a growing number of powerful open-source models (An et al., 2024; Gao et al., 2022; Radford et al., 2023). Together, these components form a virtuous cycle in which data, benchmarks, and models mutually reinforce one another, continuously driving rapid progress in both academic research and practical deployment. In contrast, the Wu dialect ², a crucial Chinese dialect that remains relatively underexplored in speech technology community, is associated with a severely underdeveloped speech processing ecosystem. The lack of sufficient datasets, standardized benchmarks, and accessible models has substantially hindered both research advancement and real-world adoption. To bridge this gap and promote inclusive speech technology, this work lays the foundation for a comprehensive speech processing ecosystem for the Wu dialect.

The Wu dialect is a major branch of Chinese spoken by approximately 100 million speakers across Shanghai, Zhejiang, Jiangsu, and overseas communities, and is linguistically characterized by exceptional complexity. It preserves the fully voiced consonant system inherited from Old Chinese. It features a highly intricate tone sandhi system in

¹<https://anonymous.4open.science/r/WenetSpeech-Wu>

²https://en.wikipedia.org/wiki/Wu_Chinese

Table 1: Comparison of typical low-resource speech processing resources related to WenetSpeech-Wu. ✓, △, and ✗ indicate experimentally verified availability, unverified availability, and absence of dimension, respectively. Abbreviations: Data Comp (Data Composition; CN = China, SEA = Southeast Asia), Hrs (hours), Trans (transcription), Para (paralinguistics), Transl (translation), Q Tier (quality grading tier), Instr TTS (instruct TTS), Sp Und. (speech understanding), Wu-SH (Shanghai sub-dialect of Wu).

Resource	Dataset						Benchmark					Model			
	Data Comp	Hrs	Trans	Para	Transl	Q Tier	ASR	AST	Para	TTS	Instr TTS	ASR	TTS	Sp Und	Instr TTS
KeSpeech	8 CN Accents	1.5k	✓	△	✗	✗	✓	✗	✗	✗	✗	✓	✗	✗	✗
GigaSpeech2	3 SEA Langs	30k	✓	✗	✗	✗	✓	✗	✗	✗	✗	✓	✗	✗	✗
WenetSpeech-Yue	CN Yue	21k	✓	△	✗	✗	✓	✗	✗	✓	✗	✓	✓	✗	✗
WenetSpeech-Chuan	CN Sichuan	10k	✓	△	✗	✗	✓	✗	✗	✓	✗	✓	✓	✗	✗
Magicdata-Shanghai	CN Wu-SH	4	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
WenetSpeech-Wu	CN Wu	8k	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

connected speech (Rose, 2015), both of which pose substantial challenges for speech modeling. Moreover, substantial variation exists among its sub-dialects, such as Shanghainese, Suzhounese, and Hangzhounese, further complicating the development of robust speech processing systems.

Despite its linguistic significance and broad speaker base, the Wu dialect remains severely under-resourced for speech processing. From a data perspective, existing open-source resources are extremely limited in both scale and coverage: the only publicly available dataset, MagicData-Shanghai³, provides merely 4.19 hours of annotated Shanghainese speech for automatic speech recognition (ASR), with no open datasets available for other Wu sub-dialects. Moreover, high-quality emotion annotations, speaker attributes and expressive speech-text pairs are entirely absent, which are essential for supporting a wide range of speech processing tasks. From an evaluation perspective, the lack of publicly available benchmarks prevents fair comparison and systematic assessment across different methods. At the model level, even foundational speech processing systems, such as ASR and text-to-speech synthesis (TTS) models, are either unavailable or perform poorly, making both open-source and commercial systems essentially unusable for Wu dialect applications. Collectively, these limitations underscore an urgent need to establish a comprehensive Wu dialect speech processing ecosystem, so as to facilitate dialectal speech research and foster the growth of the open-source community.

To address these challenges, we present a com-

prehensive open-source effort that lays the foundation for a Wu dialect speech processing ecosystem. Specifically, we construct **WenetSpeech-Wu**, the first large-scale, multi-dimensionally annotated open-source corpus for the Wu dialect, comprising approximately 8,000 hours of eight Wu sub-dialects. Building upon this dataset, we introduce **WenetSpeech-Wu-Bench**, the first standardized and publicly accessible benchmark designed for systematic evaluation of Wu dialect speech processing, covering ASR, Wu-to-Mandarin automatic speech translation (AST), speaker attribute prediction (gender and age), speech emotion recognition, TTS, and instruction-following TTS (instruct TTS). Furthermore, we release a suite of strong open-source models trained on WenetSpeech-Wu, including ASR, TTS, unified speech understanding, and instruct TTS models, which substantially outperform existing open-source and commercial systems, thereby establishing strong models and empirically demonstrating the effectiveness of the proposed dataset. Together, these contributions lay the foundation for a comprehensive, publicly accessible ecosystem for Wu dialect speech processing.

Our main contributions are threefold:

- We release **WenetSpeech-Wu**, the first large-scale and multi-dimensionally annotated open-source corpus for the Wu dialect, containing approximately 8,000 hours of speech data.
- We introduce **WenetSpeech-Wu-Bench**, the first standardized benchmark for evaluating Wu dialect speech processing, covering ASR, Wu-to-Mandarin AST, speaker attribute prediction, speech emotion recognition, TTS, and instruct TTS tasks.

³<https://magichub.com/datasets/shanghai-dialect-conversational-speech-corpus/>

- We open-source strong models trained on WenetSpeech-Wu, demonstrating substantial improvements over existing systems and validating the effectiveness of proposed dataset.

2 Related Work

2.1 Speech Datasets for Low-Resource Languages and Chinese Dialects

Recently, significant efforts have been devoted to constructing speech corpora for low-resource languages and dialects, aiming to mitigate data scarcity and promote inclusive speech technologies (Shao et al., 2025). Representative large-scale corpora such as Common Voice (Ardila et al., 2020), GigaSpeech2 (Yang et al., 2024), and KeSpeech (Tang et al., 2021) have substantially advanced research on speech processing for low-resource or underrepresented languages. Within the Chinese dialect family, several dialect-oriented speech resources have been released recently, including WenetSpeech-Yue (Li et al., 2025) for Cantonese and WenetSpeech-Chuan (Dai et al., 2025) for Sichuanese. However, these datasets primarily focus on ASR and, to a limited extent, TTS: Common Voice and GigaSpeech2 provide only transcriptions, KeSpeech augments transcriptions with dialect and basic speaker labels, and WenetSpeech-Yue and WenetSpeech-Chuan include paralinguistic annotations but evaluate their resources only through ASR and TTS experiments, without validating the effectiveness of the paralinguistic labels. As a result, the usability of such resources for broader speech processing tasks, such as speech emotion recognition, AST, and instruct TTS, remains limited.

Despite its linguistic and practical importance, the Wu dialect remains severely under-resourced in terms of publicly available speech data. To the best of our knowledge, MagicData-Shanghai is the only open-source Wu dialect dataset, providing merely 4.19 hours of annotated Shanghainese speech for ASR. Such extreme data scarcity, together with the lack of diverse annotations has substantially hindered progress in Wu dialect speech processing research.

2.2 Wu Dialect Speech Processing: Models and Evaluation

Early studies on Wu dialect speech processing mainly focused on ASR or TTS for individual sub-dialects, most notably Shanghainese, while largely

ignoring other major variants such as Suzhounese and Hangzhounese. More recent works incorporate the Wu dialect as a minor component within large-scale multilingual ASR or TTS systems, including Dolphin (Meng et al., 2025), Qwen3-ASR, Qwen3-TTS, Step-Audio2 (Wu et al., 2025), and Qwen3-Omni (Xu et al., 2025a). In these systems, Wu dialect speech is treated as a low-priority subset rather than a primary research target, resulting in foundational ASR and TTS models whose performance remains insufficient for practical Wu dialect applications. Beyond ASR and TTS, speech understanding models and instructing TTS systems tailored to the Wu dialect remain largely unexplored.

More importantly, even for basic ASR and TTS tasks, there is currently no publicly available and standardized benchmark for evaluation. As a consequence, existing studies typically evaluate models on private test sets, severely hindering fair comparison, reproducibility, and systematic progress.

In summary, as illustrated in Tabel 1, these limitations motivate the introduction of WenetSpeech-Wu, which aims to establish a complete Wu dialect speech processing ecosystem by jointly advancing datasets, benchmarks, and models.

3 WenetSpeech-Wu

3.1 Data Construction Pipeline

We propose an automatic and scalable pipeline for constructing a large-scale Wu dialect speech dataset with multi-dimensional annotations, as illustrated in Figure 2. The pipeline is designed to enable efficient data collection, robust automatic transcription, and diverse downstream annotations.

Data Collection and Filtering. We collect large-scale in-the-wild Wu dialect speech from diverse domains and sub-dialects. Non-Wu data are first removed based on metadata filtering, followed by WebRTC VAD-based segmentation. We further apply quality filtering using DNSMOS and signal-to-noise ratio (SNR), resulting in a high-quality speech corpus.

Annotation Tool Construction. To support large-scale automatic transcription, we fine-tune two pretrained ASR models using 880 hours of manually annotated Wu dialect speech. Specifically, we fine-tune Tele-CTC-FT⁴, a Connectionist Temporal Classification (CTC)-based dialect self-supervised learning (SSL) model, and Step-

⁴<https://huggingface.co/Tele-AI/TeleSpeech-ASR1.0>

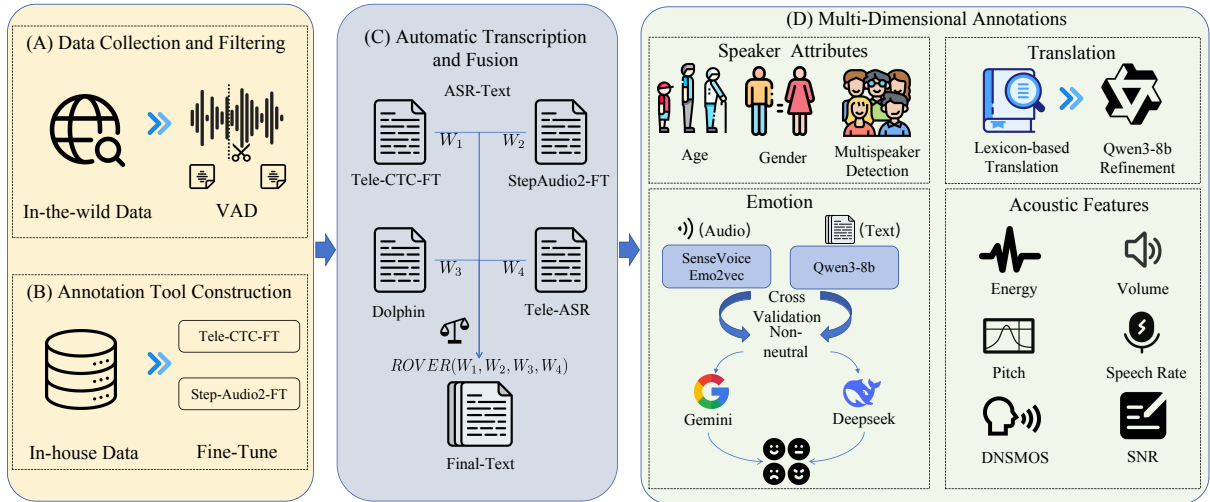


Figure 2: Data construction pipeline for WenetSpeech-Wu.

237 Audio2-FT (Wu et al., 2025). These models are
 238 used as complementary automatic annotators for
 239 Wu dialect transcription.

240 **Automatic Transcription and Fusion.** We
 241 adopt Recognizer Output Voting Error Reduction
 242 (ROVER) (Fiscus, 1997) to fuse transcription from
 243 multiple ASR systems. Specifically, we combine
 244 the outputs from the two fine-tuned Wu dialect
 245 ASR models, Dolphin (Meng et al., 2025) and
 246 TeleASR (Chen et al., 2024), and determine the
 247 model weights via grid search. The fused results
 248 provide final transcriptions along with confidence
 249 scores.

250 **Multi-Dimensional Annotations.** Each speech
 251 sample is annotated across multiple dimensions.
 252 Speaker attributes, including gender and age, are
 253 inferred using the VoxProfile (Feng et al., 2025),
 254 while multi-speaker presence is detected using
 255 Pyannote (Bredin, 2023). Wu-to-Mandarin
 256 translations are generated through lexicon-based
 257 mapping and further refined using Qwen3-8B
 258 (Yang et al., 2025), a large language model, aiming
 259 to provide fluent standard Mandarin references. Emotion
 260 annotations are obtained through a multi-stage,
 261 cross-modal procedure. Initial predictions are
 262 produced using SenseVoice (An et al., 2024) and
 263 Emo2Vec (Wang et al., 2020) for acoustic signals,
 264 and Qwen3-8B for textual content. Samples jointly
 265 identified as non-neutral are further analyzed
 266 using DeepSeek-R1 based on text and Gemini-2.5-
 267 Pro based on acoustic information, with the final
 268 label determined by the intersection of the two.
 269 In addition, prosodic acoustic features, including
 270 speech rate, loudness, energy, and pitch, are ex-

271 tracted by Dataspeech (Lyth and King, 2024) to
 272 support speech generation tasks.

3.2 Datasets 273

274 WenetSpeech-Wu is the first large-scale Wu dialect
 275 speech corpus with multidimensional annotations.
 276 It contains rich metadata and annotations, includ-
 277 ing transcriptions with confidence scores, Wu-to-
 278 Mandarin translations, domain and sub-dialect la-
 279 bels, speaker attributes, emotion annotations, and
 280 audio quality measures. The dataset comprises ap-
 281 proximately 8,000 hours of speech collected from
 282 diverse domains and covers eight Wu sub-dialects.
 283 To support a wide range of speech processing tasks
 284 with heterogeneous quality requirements, we fur-
 285 ther adopt a task-specific data quality grading strat-
 286 egy. In the following, we describe the dataset statis-
 287 tics, domain and dialect coverage, annotation distri-
 288 butions, and quality control mechanisms in detail.

289 **Duration and Confidence Distribution.**
 290 WenetSpeech-Wu contains 8,000 hours of speech
 291 with 3.86M utterances, with utterance durations
 292 up to 30 seconds and an average duration of 7.45
 293 seconds. We use the transcription confidence
 294 produced by the weighted ROVER as a measure
 295 of annotation quality, and retain utterances
 296 with confidence scores above 0.55. The detailed
 297 distributions of utterance duration and transcription
 298 confidence are shown in Figure 3b and Figure 3c.

299 **Domain and Sub-Dialect Coverage.**
 300 WenetSpeech-Wu covers a wide range of
 301 speech domains and Wu sub-dialects. The domains
 302 include *News, Culture, Vlog, Entertainment,*
 303 *Education, Podcast, Commentary, Interview, Radio*

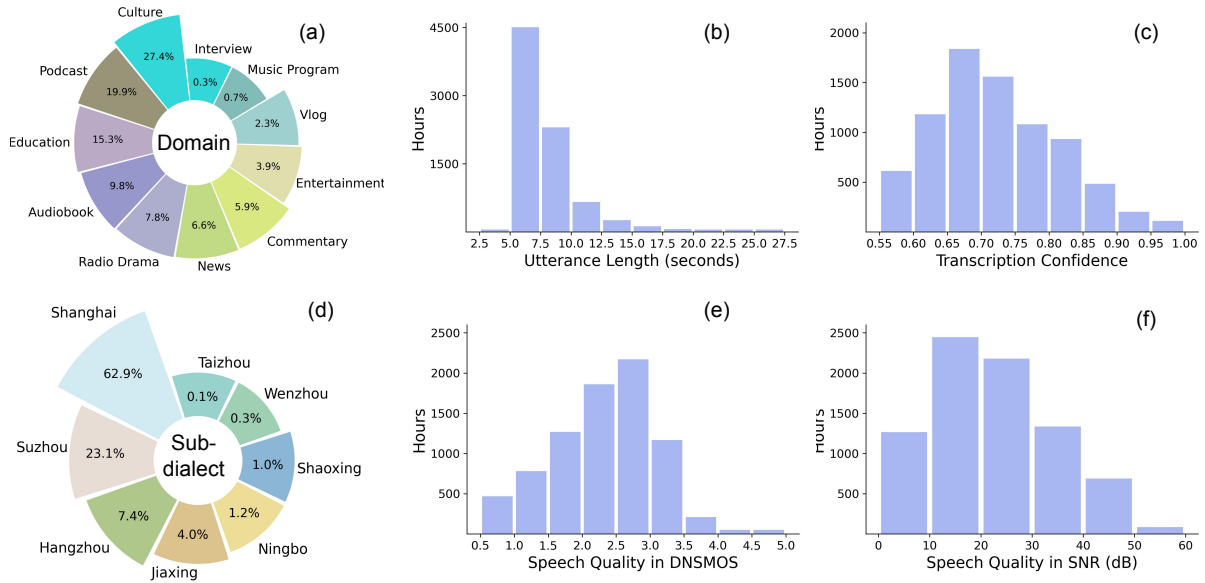


Figure 3: Statistical overview of WenetSpeech-Wu.

Table 2: Statistics of speaker attributes and emotion annotations in WenetSpeech-Wu.

Task	Category	Hours
Gender	Male	4,135
	Female	1,331
Age	Teenagers	372
	Youth	1,673
	Middle-aged	2,003
	Elderly	1,418
Emotion	Neutral	5,102
	Happy	73
	Sad	81
	Angry	109
	Surprised	101

Drama, Music Program, and Audiobook, with their distribution shown in Figure 3a. In terms of dialectal coverage, approximately 37% of the recordings cannot be reliably assigned to a specific Wu sub-dialect and are therefore labeled as Unknown. The remaining recordings span multiple identified Wu sub-dialects, including *Shanghainese, Suzhounese, Shaoxingnese, Ningbonese, Hangzhounese, Jiaxingnese, Taizhounese, and Wenzhounese*, whose distribution is illustrated in Figure 3d.

Audio Quality. As shown in Figure 3e and Figure 3f, most utterances have SNRs between 10 and 40 dB, with a peak at 20–30 dB. MOS scores are primarily in the range 2.0–3.5.

Speaker Attributes and Emotion Annotations. We annotate gender, age, and emotion labels for

single-speaker segments. Gender is categorized into *Male* and *Female*, age into four groups including 0–17 as *Teenagers*, 18–35 as *Youth*, 36–59 as *Middle-aged*, and 60+ as *Elderly*, and emotion into five classes: *Neutral, Happy, Sad, Surprised, and Angry*. The distribution of each category is shown in Table 2.

Task-Specific Data Quality Grading. To support practical training across heterogeneous speech tasks, we adopt a task-specific data-quality grading strategy aligned with the task-specific quality requirements. For ASR and TTS, we construct two quality tiers. The normal-quality subset is designed for large-scale pretraining and prioritizes data coverage and diversity, requiring only moderate transcription confidence. In contrast, the high-quality subset targets supervised fine-tuning (SFT) and applies stricter filtering criteria, including high transcription confidence, clean acoustic conditions, and reliable speaker segmentation, to provide strong and stable supervision. For tasks that are more sensitive to annotation noise and semantic ambiguity, including Wu-to-Mandarin AST, speaker attribute prediction, speech emotion recognition, TTS, and instruct TTS, we adopt stringent selection criteria, such as single-speaker recordings, high MOS, SNR, and pitch standard deviation, and verified annotation consistency, as shown in Table 3.

4 WenetSpeech-Wu-Bench

We introduce **WenetSpeech-Wu-Bench**, the first publicly available, manually curated benchmark

Table 3: Task-specific data selection and quality tiers in WenetSpeech-Wu. Abbreviations: Q Tier (quality grading tier), Text Conf (transcription confidence), Spk (speaker type, Mono = single-speaker), Pitch Std. (pitch standard deviation), Expr (expressive data obtained via emotion label), Spk Attr (speaker attributes), Emo (emotion), Inst Pro (instruct TTS prosodic control), Inst Emo (instruct TTS emotional control).

Task	Hrs	Q Tier	Text Conf	Spk	SNR	Pitch Std.	Expr
ASR	7388	Mid	> 0.6	-	-	-	-
	795	High	> 0.85	-	-	-	-
AST	795	High	> 0.85	-	-	-	-
Spk Attr	2986	High	-	Mono	-	-	-
Emo	500	High	-	Mono	>10	>50	✓
	7388	Mid	> 0.6	-	-	-	-
TTS	1500	High	> 0.65	Mono	>10	>50	-
	Inst Pro	679	High	> 0.7	Mono	>30	-
Inst Emo	161	High	> 0.7	Mono	>10	>50	✓

for Wu dialect speech processing, covering ASR, Wu-to-Mandarin AST, speaker attributes, emotion recognition, TTS, and instruct TTS, and providing a unified platform for fair evaluation.

Automatic Speech Recognition. The ASR test set of WenetSpeech-Wu-Bench comprises 9.75 hours of Shanghainese, Suzhounese, and Mandarin code-mixed speech, including single-speaker and multi-speaker scenarios. Performance is evaluated using character error rate (CER), with character-level errors reported for detailed analysis.

Wu-to-Mandarin Speech Translation. We construct a Wu-to-Mandarin AST test set. The test set contains 3000 Wu dialect utterances totaling 4.4 hours with manually verified standard Mandarin translations, covering multiple domains. Translation quality is evaluated using the BLEU score ⁵.

Speaker Attribute Prediction and Speech Emotion Recognition. This test set evaluates the prediction of age, gender, and emotion for Wu dialect speech. For speaker attributes, gender is coded as male or female, with 1,500 samples per class. And age into four groups: teenagers aged 17 and under, youth aged 18 to 35, middle-aged adults aged 36 to 59, and elderly aged 60 and above, with 500 samples in each group. For emotion, there are 300 neutral samples, 200 happy samples, 100

sad samples, 200 angry samples, and 200 surprised samples, for a total of 1,000 samples. Performance is measured by category-wise and overall classification accuracy.

Text to Speech. We constructed a TTS testset as part of WenetSpeech-Wu-Bench, comprising 144 easy and 98 hard test sentences. The texts were further reviewed and refined by professional Wu dialect experts. Prompt audio samples were selected from the open-source Magicdata-Shanghai, and 12 speakers of the Wu dialect were selected using strict filtering criteria. For evaluation, speaker similarity is measured using WeSpeaker ⁶ based speaker embedding similarity, while intelligibility is assessed by computing CER with our proposed Step-Audio2-Wu-ASR model. Additionally, subjective listening tests are conducted, including intelligibility MOS (IMOS), similarity MOS (SMOS), and accent MOS (AMOS). The subjective evaluation was conducted with 23 listeners, each rating 20 selected samples.

Instruct TTS. WenetSpeech-Wu-Bench includes two evaluation sets for instruct TTS. The prosodic control test set consists of five speech prompts, each spoken at a moderate speaking rate and normal fundamental frequency, and synthesized into 20 sentences with controlled variations in speaking rate and pitch. For evaluation, two experiments are conducted: one with a fast speaking rate and high pitch, and another with a slow rate and low pitch. These samples are automatically annotated using Dataspeech (Lyth and King, 2024). Each pair is scored as one if the relative speaking rate and pitch match the intended instructions, and zero otherwise. The prosodic classification metric is the average score across all pairs, reflecting the model’s ability to follow prosodic instructions (Kuan et al., 2023). The emotional control test set evaluates a model’s ability to follow emotion-related instructions. We select 10 reference prompts that contain no explicit emotional expression. Based on each prompt, 50 sentences are synthesized for each of the four target emotions: anger, sadness, happiness, and surprise. The samples are evaluated using our Step-Audio2-Wu-Und model. A sample is considered correct if the predicted emotion matches the intended target emotion, and the mean classification accuracy is reported as the evaluation metric (Gao et al., 2025). Additionally, subjective listening tests are conducted to assess the quality of instruction-

⁵<https://github.com/mjpost/sacrebleu>

⁶<https://github.com/wenet-e2e/wespeaker>

Table 4: ASR results (CER%) on various test sets. Gray, red, light green, and dark green rows denote open-source baselines, commercial models, ASR models trained on WenetSpeech-Wu, and annotation models trained on in-house data. **Bold** numbers indicate best results; underlined numbers indicate second-best results.

Model	In-House		WS-Wu-Bench
	Dialogue	Reading	ASR
ASR Models			
Paraformer	63.13	66.85	64.92
SenseVoice-small	29.20	31.00	46.85
Whisper-medium	79.31	83.94	78.24
FireRedASR-AED-L	51.34	59.92	56.69
Step-Audio2-mini	24.27	24.01	26.72
Qwen3-ASR	23.96	24.13	29.31
Tencent-Cloud-ASR	23.25	25.26	29.48
Gemini-2.5-pro	85.50	84.67	89.99
Conformer-U2pp-Wu	15.20	12.24	15.14
Whisper-medium-Wu	14.19	11.09	<u>14.33</u>
Step-Audio2-Wu-ASR	<u>8.68</u>	7.86	12.85
Annotation Models			
Dolphin-small	24.78	27.29	26.93
TeleASR	29.07	21.18	30.81
Step-Audio2-FT	8.02	6.14	15.64
Tele-CTC-FT	11.90	<u>7.23</u>	23.85

following. Listeners rate the samples using two scores: the prosodic MOS (PMOS) (Chan and Kuang, 2025) and emotional Mean Opinion Score (EMOS) (Cho et al., 2025). The evaluation involves 23 listeners, each rating 15 samples, to provide perceptual judgments on how well the synthesized speech follows the intended emotional instructions.

5 Models & Experiments

To address the lack of Wu dialect speech processing models, we develop models for the two core aspects of speech processing, speech understanding and speech generation trained on WenetSpeech-Wu. These include ASR models and unified speech understanding models for understanding, as well as TTS and instruct TTS models for generation.

5.1 Speech Understanding

ASR models. To accommodate different application scenarios, we develop three Wu dialect ASR models at three scales: a Conformer-U2pp-Wu (Wu et al., 2021) model, a Whisper-Medium-Wu (Radford et al., 2023) model, and a Step-Audio2-Wu-ASR model. The Conformer-U2pp-Wu and Whisper-Medium-Wu models are trained within the Wenet framework (Yao et al., 2021), with Conformer-U2pp-Wu trained from scratch, and Whisper-Medium-Wu fine-tuned from publicly

Table 5: Speech understanding performance on WenetSpeech-Wu-Bench.

Model	ASR	AST	Gender	Age	Emotion
Qwen3-Omni	44.27	33.31	0.977	0.541	0.667
Step-Audio2-mini	<u>26.72</u>	<u>37.81</u>	0.855	0.370	0.460
Step-Audio2-Wu-Und	13.23	53.13	<u>0.956</u>	0.729	0.712

available pre-trained weights. The Step-Audio2-Wu-ASR model is a Step-Audio2-mini fine-tuned in a parameter-efficient manner using LoRA (Hu et al., 2021) within the MS-Swift framework (Zhao et al., 2025). All models are pretrained on the ASR-Mid subset and conduct SFT on the ASR-High subset. Evaluation is performed on the ASR test set of WenetSpeech-Wu-Bench as well as two in-house manually annotated test sets covering dialogue and reading scenarios, enabling comprehensive assessment across diverse speaking conditions.

For comparison, we include both open-source and commercial baselines, such as Dolphin, SenseVoice, Paraformer (Gao et al., 2022), FireRedASR (Xu et al., 2025b) as an open-source baseline, Qwen3-ASR, Gemini-2.5-Pro, TeleASR (Chen et al., 2024), and Tencent-Cloud-ASR as a commercial baseline. The two in-house test sets are strictly out-of-set evaluations for Conformer-U2pp-Wu, Whisper-medium-Wu, and Step-Audio2-Wu-ASR. Conversely, the WS-Wu-Bench-ASR test set serves as an out-of-set evaluation for Step-Audio2-FT and Tele-CTC-FT, which are trained exclusively on in-house data.

As shown in Table 4, all existing open-source and commercial ASR systems perform poorly across all three test sets, indicating that they are not viable for Wu dialect recognition. In contrast, models trained on WenetSpeech-Wu achieve state-of-the-art performance across all model scales, with even the smallest Conformer-U2pp-Wu substantially outperforming all prior systems.

Speech understanding models. For speech understanding, we train the Step-Audio2-Wu-Und model on the WenetSpeech-Wu corpus using task-aware, quality-graded data. The model is first pre-trained on ASR-Mid and High-quality AST and age, gender, and emotion annotations, and then fine-tuned with the ASR-high subset alongside the same non-ASR tasks.

For benchmarking, we compare against baseline models, including Step-Audio2-mini and Qwen3-Omni, for speech understanding tasks. As shown in Table 5, after multi-task fine-tuning, the ASR per-

Table 6: TTS results on WenetSpeech-Wu-Bench. **Bold** and underlined values denote the best and second-best results, respectively; light green rows indicate models trained on WenetSpeech-Wu or further fine-tuned on an internal high-quality dataset.

Model	WS-Wu-Eval-TTS-easy					WS-Wu-Eval-TTS-hard				
	CER(%)↓	SIM↑	IMOS↑	SMOS↑	AMOS↑	CER(%)↓	SIM↑	IMOS↑	SMOS↑	AMOS↑
Qwen3-TTS†	<u>5.95</u>	–	<u>4.35</u>	–	<u>4.19</u>	<u>16.45</u>	–	<u>4.03</u>	–	3.91
DiaMoE-TTS	57.05	0.702	3.11	3.43	3.52	82.52	0.587	2.83	3.14	3.22
CosyVoice2	10.33	0.713	3.83	3.71	3.84	82.49	<u>0.618</u>	3.24	3.42	3.37
CosyVoice2-Wu-CPT	6.35	0.727	4.01	3.84	3.92	32.97	0.620	3.72	3.55	3.63
CosyVoice2-Wu-SFT	6.19	<u>0.726</u>	4.32	<u>3.78</u>	4.11	25.00	0.601	3.96	<u>3.48</u>	3.76
CosyVoice2-Wu-SS*	5.42	–	4.37	–	4.21	15.45	–	4.04	–	<u>3.88</u>

† Commercial system with a single fixed speaker, and speaker similarity is not considered.

* Single-speaker finetuned model, and speaker similarity is not evaluated.

499 performance of Step-Audio2-Wu-Und shows a slight
500 drop than Step-Audio2-Wu-ASR but still achieves
501 the second-best results. For the Wu-to-Mandarin
502 AST task, it significantly outperforms baseline
503 models. Comparison with Step-Audio2-mini il-
504 lustrates the mismatch between Wu dialect and
505 Mandarin in gender, age, and emotion prediction,
506 which is effectively addressed by our data. Com-
507 pared with Qwen3-Omni, our model shows notable
508 improvements in age and emotion prediction, while
509 performing slightly lower on gender classification.

5.2 Speech Generation

510 **TTS models.** In the continual pre-training (CPT)
511 stage, we continue training from the open-source
512 CosyVoice2 (Du et al., 2024) checkpoint using the
513 TTS-Mid dataset for ten epochs. In the SFT stage,
514 we use TTS-High dataset trained for three epochs.
515 Finally, we perform speaker-specific fine-tuning on
516 a 10-hour internal high-quality dataset.

517 As shown in Table 6, the experimental results
518 demonstrate that the staged training strategy sig-
519 nificantly improves CosyVoice2’s speech synthe-
520 sis performance in terms of Wu dialect capabil-
521 ity. The CPT stage, leveraging large-scale pipeline-
522 processed data, enhances the model’s fundamental
523 capabilities and robustness, particularly on chal-
524 lenging samples. The SFT stage, further im-
525 proves speech naturalness and expressiveness. Fi-
526 nally, the Single-Speaker Supervised Fine-Tuning
527 (SS-SFT) stage achieves the best performance
528 across CER, IMOS, and AMOS metrics. Overall,
529 CosyVoice2-Wu-SS approaches or surpasses the
530 baseline Qwen3-TTS, DiaMoE-TTS (Chen et al.,
531 2025), and CosyVoice2 in most metrics, especially
532 in difficult sample synthesis.

533 **Instruct TTS models.** The instruction training

Table 7: Performance of instruct TTS model.

Type	Metric	CosyVoice2-Wu -SFT	CosyVoice2-Wu -instruct
Emotion	Happy↑	0.87	0.94
	Angry↑	0.83	0.87
	Sad↑	0.84	0.88
	Surprised↑	0.67	0.73
	EMOS↑	3.66	3.83
Prosody	Pitch↑	0.24	0.74
	Speech Rate↑	0.26	0.82
	PMOS↑	2.13	3.68

534 data are from the Inst Pro and Inst Emo datasets as
535 introduced in Table 3. Training was performed
536 on the instruction text and tags using a small
537 learning rate for five epochs. Compared with the
538 model before instruction fine-tuning, the results
539 showed clear improvements across all controllabil-
540 ity metrics on WenetSpeech-Wu-Bench, as shown
541 in Table 7, and subjective listening tests confirmed
542 strong perceptual control effects, validating the ef-
543 fectiveness of proposed data.

6 Conclusion

544 In this work, we establish a Wu dialect speech pro-
545 cessing ecosystem encompassing datasets, bench-
546 marks, and models. We introduce WenetSpeech-
547 Wu, an 8,000-hour large-scale corpus with rich
548 multi-dimensional annotations across multiple Wu
549 sub-dialects, and present the first public Wu dialect
550 benchmark covering ASR, Wu-to-Mandarin AST,
551 speech attribute and emotion recognition, TTS, and
552 instruct TTS. Building on this dataset, we release a
553 suite of ASR, TTS, unified speech understanding
554 models, and instruct TTS models, enabling commu-
555 nity research and demonstrating the effectiveness
556 of WenetSpeech-Wu.

559 **Limitations**

560 Despite the contributions of WenetSpeech-Wu to
561 the dialectal speech processing community, several
562 limitations remain to be addressed in future work.
563 Although WenetSpeech-Wu is large-scale and cov-
564 ers multiple Wu sub-dialects, the data distribution
565 across dialects and domains is not perfectly bal-
566 anced, which may affect model generalization to
567 less-represented varieties. In addition, many anno-
568 tations are produced through automated or semi-
569 automated pipelines. While we apply stringent
570 quality control and filtering strategies, these anno-
571 tations may still contain noise compared to fully
572 manual labeling. Finally, our baseline models are
573 intended to provide strong and reproducible ref-
574 erence points rather than task-optimal solutions,
575 and future work may further improve performance
576 through more specialized modeling and training
577 strategies.

578 **Ethical Considerations**

579 The Wu dialect speech data used in this work are
580 collected exclusively from publicly accessible, non-
581 restricted online sources intended for open dissem-
582 ination. No access control mechanisms were cir-
583 cumvented during data collection. To ensure legal
584 and ethical compliance, we only include audio that
585 is explicitly available for research use.

586 To further mitigate potential ethical and privacy
587 risks, we collaborate with a professional data ser-
588 vice company to conduct independent review and
589 filtering of the collected data. This process fo-
590 cuses on identifying and removing content that may
591 contain personally identifiable information, sensi-
592 tive attributes, or offensive material. In particular,
593 audio segments involving explicit personal identi-
594 fiers, such as names, phone numbers, or addresses,
595 are excluded. As a result, the released dataset is
596 anonymized and does not enable the identification
597 of individual speakers.

598 All benchmark annotations are produced by the
599 same professional data service company through
600 manual labeling, following established internal
601 compliance and quality control procedures. No
602 personally identifiable information is disclosed to
603 annotators, and annotation tasks are strictly limited
604 to speech content relevant to linguistic analysis and
605 benchmarking. Due to commercial confidentiality,
606 detailed information regarding annotator compen-
607 sation and internal annotation protocols cannot be
608 publicly disclosed; however, the service provider

609 confirms that all annotation activities comply with
610 applicable local labor regulations and professional
611 annotation standards.

612 The dataset and benchmark are constructed
613 solely to support academic research on Wu di-
614 alect speech processing. Based on our data col-
615 lection and filtering procedures, we think that the
616 intended use of the dataset and benchmark does
617 not introduce ethical or societal risks. To promote
618 transparency and responsible research practices, we
619 provide detailed documentation describing dataset
620 composition, scale, splits, and evaluation protocols,
621 and we clearly report experimental settings and
622 results. The dataset is released with explicit docu-
623 mentation specifying its intended non-commercial
624 research use.

625 **References**

- 626 Keyu An, Qian Chen, Chong Deng, Zhihao Du,
627 Changfeng Gao, Zhifu Gao, Yue Gu, Ting He,
628 Hangrui Hu, Kai Hu, Shengpeng Ji, Yabin Li, Zerui
629 Li, Heng Lu, Haoneng Luo, Xiang Lv, Bin Ma,
630 Ziyang Ma, Chongjia Ni, and 14 others. 2024. Funau-
631 diollm: Voice understanding and generation founda-
632 tion models for natural interaction between humans
633 and llms. *CoRR*, abs/2407.04051.
- 634 Rosana Ardila, Megan Branson, Kelly Davis, Michael
635 Kohler, Josh Meyer, Michael Henretty, Reuben
636 Morais, Lindsay Saunders, Francis M. Tyers, and
637 Gregor Weber. 2020. Common voice: A massively-
638 multilingual speech corpus. In *Proc. LREC*, pages
639 4218–4222.
- 640 Hervé Bredin. 2023. pyannote.audio 2.1 speaker di-
641 arization pipeline: principle, benchmark, and recipe.
642 In *Proc. INTERSPEECH*, pages 1983–1987.
- 643 Cedric Chan and Jianjing Kuang. 2025. Toward Ob-
644 jective and Interpretable Prosody Evaluation in Text-
645 to-Speech: A Linguistically Motivated Approach.
646 *arXiv*, 2511.02104.
- 647 Hongjie Chen, Zehan Li, Guangmin Xia, Boqing Liu,
648 Yan Yang, Jian Kang, and Jie Li. 2024. *Tele-
649 SpeechPT: Large-Scale Chinese Multi-dialect and
650 Multi-accent Speech Pre-training*, pages 183–190.
651 Springer.
- 652 Ziqi Chen, Gongyu Chen, Yihua Wang, Chaofan Ding,
653 Zihao Chen, and Wei-Qiang Zhang. 2025. Diamoe-
654 tts: A unified ipa-based dialect TTS framework with
655 mixture-of-experts and parameter-efficient zero-shot
656 adaptation. *CoRR*, abs/2509.22727.
- 657 Deok-Hyeon Cho, Hyung-Seok Oh, Seung-Bin Kim,
658 and Seong-Whan Lee. 2025. Emosphere++:
659 Emotion-controllable zero-shot text-to-speech via
660 emotion-adaptive spherical vector. *IEEE Trans. Af-
661 fect. Comput.*, 16(3):2365–2380.

662	Yuhang Dai, Ziyu Zhang, Shuai Wang, Longhao Li,	Yangyang Meng, Jinpeng Li, Guodong Lin, Yu Pu,	719
663	Zhao Guo, Tianlun Zuo, Shuiyuan Wang, Hongfei	Guanbo Wang, Hu Du, Zhiming Shao, Yukai Huang,	720
664	Xue, Chengyou Wang, Qing Wang, Xin Xu, Hui Bu,	Ke Li, and Wei-Qiang Zhang. 2025. Dolphin: A	721
665	Jie Li, Jian Kang, Binbin Zhang, and Lei Xie. 2025.	large-scale automatic speech recognition model for	722
666	Wenetspeech-chuan: A large-scale sichuanese corpus	eastern languages. <i>CoRR</i> , abs/2503.20212.	723
667	with rich annotation for dialectal speech processing.		
668	<i>CoRR</i> , abs/2509.18004.		
669	Zhihao Du, Yuxuan Wang, Qian Chen, Xian Shi, Xiang	Jing Peng, Yucheng Wang, Bohan Li, Yiwei Guo, Han-	724
670	Lv, Tianyu Zhao, Zhifu Gao, Yexin Yang, Changfeng	kun Wang, YanGui Fang, Yu Xi, Haoyu Li, Xu Li,	725
671	Gao, Hui Wang, Fan Yu, Huadai Liu, Zhengyan	Ke Zhang, Shuai Wang, and Kai Yu. 2025. A Sur-	726
672	Sheng, Yue Gu, Chong Deng, Wen Wang, Shil-	vey on Speech Large Language Models for Under-	727
673	liang Zhang, Zhijie Yan, and Jingren Zhou. 2024.	standing. <i>IEEE Journal of Selected Topics in Signal</i>	728
674	Cosyvoice 2: Scalable streaming speech synthesis	<i>Processing</i> , page 1–32.	729
675	with large language models. <i>CoRR</i> , abs/2412.10117.		
676	Tiantian Feng, Jihwan Lee, Anfeng Xu, Yoonjeong	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brock-	730
677	Lee, Thanathai Lertpetchpun, Xuan Shi, Helin Wang,	man, Christine McLeavey, and Ilya Sutskever. 2023.	731
678	Thomas Thebaud, Laureano Moro-Velázquez, Dani	Robust speech recognition via large-scale weak su-	732
679	Byrd, Najim Dehak, and Shrikanth Narayanan. 2025.	pervision. In <i>Proc. ICML</i> , volume 202, pages 28492–	733
680	Vox-profile: A speech foundation model benchmark	28518.	734
681	for characterizing diverse speaker and speech traits.		
682	<i>CoRR</i> , abs/2505.14648.	Philip Rose. 2015. Tonation in three chinese wu dialects.	735
683	J.G. Fiscus. 1997. A post-processing system to yield	In <i>Proc. ICPhS</i> .	736
684	reduced word error rates: Recognizer output voting		
685	error reduction (rover). In <i>Proc. ASRU</i> .	S. Sakshi, Utkarsh Tyagi, Sonal Kumar, Ashish Seth,	737
686	Changfeng Gao, Zhihao Du, and Shiliang Zhang. 2025.	Ramaneswaran Selvakumar, Oriol Nieto, Ramani Du-	738
687	Differentiable reward optimization for LLM based	raiswami, Sreyan Ghosh, and Dinesh Manocha. 2025.	739
688	TTS system. In <i>Proc. Interspeech</i> .	MMAU: A massive multi-task audio understanding	740
689	Zhifu Gao, Shiliang Zhang, Ian McLoughlin, and Zhijie	and reasoning benchmark. In <i>Proc. ICLR</i> .	741
690	Yan. 2022. Paraformer: Fast and accurate parallel		
691	transformer for non-autoregressive end-to-end speech	Mingchen Shao, Bingshen Mu, Chengyou Wang,	742
692	recognition. In <i>Proc. Interspeech</i> , pages 2063–2067.	Haizhou Li, Ying Yan, Zhonghua Fu, and Lei	743
693	Haorui He, Zengqiang Shang, Chaoren Wang, Xuyuan	Xie. 2025. Towards Building Speech Large Lan-	744
694	Li, Yicheng Gu, Hua Hua, Liwei Liu, Chen Yang,	guage Models for Multitask Understanding in Low-	745
695	Jiaqi Li, Peiyang Shi, Yuancheng Wang, Kai Chen,	Resource Languages. <i>CoRR</i> , abs/2509.14804.	746
696	Pengyuan Zhang, and Zhizheng Wu. 2024. Emilia:		
697	An Extensive, Multilingual, and Diverse Speech	Zhiyuan Tang, Dong Wang, Yanguang Xu, Jianwei Sun,	747
698	Dataset For Large-Scale Speech Generation. In <i>Proc.</i>	Xiaoning Lei, Shuaijiang Zhao, Cheng Wen, Xingjun	748
699	<i>SLT</i> , pages 885–890.	Tan, Chuandong Xie, Shuran Zhou, Rui Yan, Chen-	749
700	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	jia Lv, Yang Han, Wei Zou, and Xiangang Li. 2021.	750
701	Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu	KeSpeech: An Open Source Speech Dataset of Man-	751
702	Chen. 2021. Lora: Low-rank adaptation of large	darin and Its Eight Subdialects. In <i>Proc. NeurIPS</i> .	752
703	language models. <i>CoRR</i> , abs/2106.09685.		
704	Chun-Yi Kuan, Chen-An Li, Tsu-Yuan Hsu, Tse-Yang	Dingdong Wang, Jincenzi Wu, Junan Li, Dongchao	753
705	Lin, Ho-Lam Chung, Kai-Wei Chang, Shuo-Yiin	Yang, Xueyuan Chen, Tianhua Zhang, and Helen	754
706	Chang, and Hung-Yi Lee. 2023. Towards general-	Meng. 2025a. MMSU: A massive multi-task spoken	755
707	purpose text-instruction-guided voice conversion. In	language understanding and reasoning benchmark.	756
708	<i>Proc. ASRU</i> , pages 1–8.	<i>CoRR</i> , abs/2506.04779.	757
709	Longhao Li, Zhao Guo, Hongjie Chen, Yuhang Dai,	Shuai Wang, Zhaokai Sun, Zhennan Lin, Chengyou	758
710	Ziyu Zhang, Hongfei Xue, Tianlun Zuo, Chengyou	Wang, Zhou Pan, and Lei Xie. 2025b. MSU-Bench:	759
711	Wang, Shuiyuan Wang, Jie Li, Jian Kang, Xin Xu,	Towards Understanding the Conversational Multi-	760
712	Hui Bu, Binbin Zhang, Ruibin Yuan, Ziya Zhou,	talker Scenarios. <i>CoRR</i> , abs/2508.08155.	761
713	Wei Xue, and Lei Xie. 2025. Wenetspeech-yue:		
714	A large-scale cantonese speech corpus with multi-	Shuo Wang, Aishan Maoliniazhi, Xinle Wu, and Xi-	762
715	dimensional annotation. <i>CoRR</i> , abs/2509.03959.	aofeng Meng. 2020. Emo2vec: Learning emotional	763
716	Daniel Lyth and Simon King. 2024. Natural language	embeddings via multi-emotion category. <i>ACM Trans.</i>	764
717	guidance of high-fidelity text-to-speech with syn-	<i>Internet Techn.</i> , 20(2):13:1–13:17.	765
718	thetic annotations. <i>CoRR</i> , abs/2402.01912.	Boyong Wu, Chao Yan, Chen Hu, Cheng Yi, Chengli	766
		Feng, Fei Tian, Feiyu Shen, Gang Yu, Haoyang	767
		Zhang, Jingbei Li, Mingrui Chen, Peng Liu, Wang	768
		You, Xiangyu Tony Zhang, Xingyuan Li, Xuerui	769
		Yang, Yayue Deng, Yechang Huang, Yuxin Li, and 81	770
		others. 2025. Step-audio 2 technical report. <i>CoRR</i> ,	771
		abs/2507.16632.	772

773 Di Wu, Binbin Zhang, Chao Yang, Zhendong
774 Peng, Wenjing Xia, Xiaoyu Chen, and Xin Lei.
775 2021. U2++: unified two-pass bidirectional end-
776 to-end model for speech recognition. *CoRR*,
777 abs/2106.05642.

778 Jin Xu, Zhifang Guo, Hangrui Hu, Yunfei Chu, Xiong
779 Wang, Jinzheng He, Yuxuan Wang, Xian Shi, Ting
780 He, Xinfu Zhu, Yuanjun Lv, Yongqi Wang, Dake
781 Guo, He Wang, Linhan Ma, Pei Zhang, Xinyu
782 Zhang, Hongkun Hao, Zishan Guo, and 19 others.
783 2025a. Qwen3-Omni Technical Report. *CoRR*,
784 abs/2509.17765.

785 Kaituo Xu, Feng-Long Xie, Xu Tang, and Yao Hu.
786 2025b. Fireredas: Open-source industrial-grade
787 mandarin speech recognition models from encoder-
788 decoder to LLM integration. *CoRR*, abs/2501.14350.

789 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,
790 Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,
791 Chengen Huang, Chenxu Lv, Chujie Zheng, Day-
792 iheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao
793 Ge, Haoran Wei, Huan Lin, Jialong Tang, and 40
794 others. 2025. Qwen3 technical report. *CoRR*,
795 abs/2505.09388.

796 Yifan Yang, Zheshu Song, Jianheng Zhuo, Mingyu Cui,
797 Jinpeng Li, Bo Yang, Yexing Du, Ziyang Ma, Xun-
798 ying Liu, Ziyuan Wang, Ke Li, Shuai Fan, Kai Yu, Wei-
799 Qiang Zhang, Guoguo Chen, and Xie Chen. 2024.
800 Gigaspeech 2: An evolving, large-scale and multi-
801 domain ASR corpus for low-resource languages with
802 automated crawling, transcription and refinement.
803 *CoRR*, abs/2406.11546.

804 Zhuoyuan Yao, Di Wu, Xiong Wang, Binbin Zhang,
805 Fan Yu, Chao Yang, Zhendong Peng, Xiaoyu Chen,
806 Lei Xie, and Xin Lei. 2021. Wenet: Production
807 oriented streaming and non-streaming end-to-end
808 speech recognition toolkit. In *Proc. Interspeech*,
809 pages 4054–4058.

810 Heiga Zen, Viet Dang, Rob Clark, Yu Zhang, Ron J.
811 Weiss, Ye Jia, Zhifeng Chen, and Yonghui Wu. 2019.
812 LibriTTS: A Corpus Derived from LibriSpeech for
813 Text-to-Speech. In *Proc. Interspeech*, pages 1526–
814 1530.

815 Binbin Zhang, Hang Lv, Pengcheng Guo, Qijie Shao,
816 Chao Yang, Lei Xie, Xin Xu, Hui Bu, Xiaoyu Chen,
817 Chenchen Zeng, Di Wu, and Zhendong Peng. 2022.
818 WENETSPEECH: A 10000+ Hours Multi-Domain
819 Mandarin Corpus for Speech Recognition. In *Proc.*
820 *ICASSP*, pages 6182–6186.

821 Yuze Zhao, Jintao Huang, Jinghan Hu, Xingjun Wang,
822 Yunlin Mao, Daoze Zhang, Zeyinzi Jiang, Zhikai Wu,
823 Baole Ai, Ang Wang, Wenmeng Zhou, and Yingda
824 Chen. 2025. SWIFT: A scalable lightweight infras-
825 tructure for fine-tuning. In *Proc. AAAI*, pages 29733–
826 29735.

A Appendix 827

A.1 Details of Datasets 828

Example of an Annotated Data. Each speech sample is represented using a unified JSON format. An entry includes a unique key, audio duration, Wu dialect transcription with its Mandarin translation, and an associated confidence score. Contextual metadata, such as domain and sub-dialect, audio quality indicators (e.g., DNSMOS and SNR), and paralinguistic annotations (emotion, age, and gender) are organized into structured fields. In addition, prosodic acoustic features, including speech rate, loudness, fundamental frequency statistics, and energy statistics, are provided to support fine-grained acoustic analysis. Figure 4 illustrates an example of the annotated JSON entry. 829 830 831 832 833 834 835 836 837 838 839 840 841 842

```
{
  "key": "mx9ok0979_qtmvx_8765690",
  "duration": "2.3",
  "transcription": "依好今朝天气蛮好",
  "translation": "你好今天天气挺好的",
  "confidence": "0.89",
  "meta-info": {
    "domain": "Interview",
    "sub-dialect": "Shanghai",
  },
  "audio-quality": {
    "DNSMOS": "3.98",
    "SNR": "58.74",
  },
  "paralinguistic": {
    "emotion": "Happy",
    "age": "Youth",
    "gender": "Male",
  },
  "low-level-feature": {
    "speech_rate": 15.094,
    "loudness": -17.72,
    "f0_mean": 167.17,
    "f0_std": 40.74,
    "energy_mean": 0.0684,
    "energy_std": 0.078,
  }
}
```

Figure 4: Example of an Annotated Data Entry in JSON format

A.2 Statistics of the WenetSpeech-Wu-Bench 843

As summarized in Table 9, the WenetSpeech-Wu Bench consists of diverse test sets tailored for multiple speech-related tasks. For each task, we report the number of utterances and the corresponding total duration in hours. 844 845 846 847 848

A.3 Details of Experiments 849

Optimization and Training Hyperparameters. The training hyperparameters for different models 850 851

Table 8: Optimization and training hyperparameters for proposed ASR and TTS models.

Model	Para. (M)	LR	LR sched.	Warmup	Grad. Acc.	Batch size
Whisper-Medium-Wu	769	8×10^{-5}	WarmupLR	4000 steps	4	Dynamic (24k frames)
Conformer-U2pp-Wu	123	1×10^{-3}	WarmupLR	25000 steps	4	Dynamic (60k frames)
Step-Audio2-Wu-ASR	7000	1×10^{-5}	WarmupLR	0.05 ratio	8	8
CosyVoice2-Wu-CPT	500	1×10^{-4}	WarmupLR	25000 steps	2	Dynamic (10k frames)
CosyVoice2-Wu-SFT	500	1×10^{-5}	ConstantLR	0	2	Dynamic (2k frames)
CosyVoice2-Wu-SS	500	5×10^{-6}	ConstantLR	0	2	Dynamic (1k frames)
CosyVoice2-Wu-instruct	500	5×10^{-6}	ConstantLR	0	2	Dynamic (1k frames)

Table 9: Statistics of the WenetSpeech-Wu-Bench.

Task	Category	Utterances	Duration (h)
ASR	-	4851	9.75
AST	-	3000	4.4
Gender	Male	1500	2.10
	Female	1500	2.29
	Total	3000	4.39
Age	Youth	500	0.69
	Middle-aged	500	0.79
	Elderly	500	0.69
	Total	1500	2.22
Emotion	Neutral	300	0.50
	Happy	200	0.23
	Sad	100	0.26
	Angry	200	0.15
	Surprised	200	0.27
	Total	1000	1.41

Prompt Templates Used in Experiments. For all speech understanding tasks, the same prompt templates are used during both training and inference to avoid potential distribution mismatch. For instruct TTS tasks, instruction-based templates are likewise employed in both training and inference to enable explicit instruction control. Table 10 lists the prompt formulations for each task. The English prompts are shown for ease of presentation, while equivalent instructions are employed in other languages when applicable.

Table 10: Details of prompt templates used in speech understanding and instruct TTS experiments.

Task Type	Prompt Content
ASR	Please transcribe the speech.
AST	Please listen to this speech carefully and translate its content into Mandarin.
Age	Based on the acoustic features of the speech, determine the speaker’s age. Choose one label from: youth, middle-aged, or elderly.
Gender	Based on the acoustic features of the speech, determine the speaker’s gender. Choose one label from: male or female.
Emotion	Based on the acoustic features and semantics of the speech, determine the emotion. Choose one label from: neutral, happy, sad, surprised, or angry.
Inst Emo	You must speak with anger, sadness, happiness, or surprise.
Inst Pro	This is a man or woman speaking in a low or high voice and at a slow or fast pace.

are summarized in Table 8.

Supplementary Experimental Results. To provide a more detailed analysis of model performance on WS-Wu-Bench, we report category-wise classification accuracy for gender, age, and emotion across WS-Wu-Und and baseline models. As shown in Figure 5, our model achieves more balanced and consistently higher performance across different paralinguistic categories, demonstrating stronger paralinguistic recognition capability.

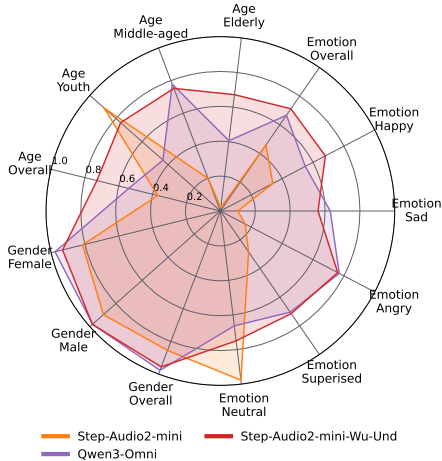


Figure 5: Supplementary comparison results on age, gender, and emotion recognition tasks.