

MRSAudio: A Large-Scale Multimodal Recorded Spatial Audio Dataset with Refined Annotations

Wenxiang Guo* Changhao Pan* Zhiyuan Zhu* Xintong Hu* Yu Zhang* Li Tang
Rui Yang Han Wang Zongbao Zhang Yuhua Wang Yixuan Chen Hankun Xu Ke Xu
Pengfei Fan Zhetao Chen Yanhao Yu Qiange Huang Fei Wu Zhou Zhao[†]
Zhejiang University
{guowx314,panch,zhaozhou}@zju.edu.cn



Figure 1: Overview of MRSAudio. The dataset comprises four real-world scenarios: MRSSpeech, MRSLife, MRSMusic, and MRSSing, each with multimodal annotations for spatial audio research.

Abstract

Humans rely on multisensory integration to perceive spatial environments, where auditory cues enable sound source localization in three-dimensional space. Despite the critical role of spatial audio in immersive technologies such as VR/AR, most existing multimodal datasets provide only monaural audio, which limits the development of spatial audio generation and understanding. To address these challenges, we introduce MRSAudio, a large-scale multimodal spatial audio dataset designed to advance research in spatial audio understanding and generation. MRSAudio spans four distinct components: MRSLife, MRSSpeech, MRSMusic, and MRSSing, covering diverse real-world scenarios. The dataset includes synchronized binaural and ambisonic audio, exocentric and egocentric video, motion trajectories, and fine-grained annotations such as transcripts, phoneme boundaries, lyrics, scores, and prompts. To demonstrate the utility and versatility of MRSAudio, we establish five foundational tasks: audio spatialization, and spatial text to speech, spatial singing voice synthesis, spatial music generation and sound event localization and detection. Results show that MRSAudio enables high-quality spatial modeling and supports a broad range of spatial audio research. Demos and dataset access are available at <https://mrsaudio.github.io>.

*Equal contribution

[†]Corresponding Author

18 **1 Introduction**

19 Humans rely on multisensory integration to perceive and interpret physical environments. With the
20 rapid growth of film, virtual reality (VR), augmented reality (AR), and gaming applications, users
21 increasingly expect not only precise audiovisual alignment but also highly immersive experiences.
22 While recent advances in deep learning have enabled realistic generation of speech, music, and sound
23 effects synchronized with text or video (Kreuk et al., 2022; Dongchao Yang, 2023; Wang et al., 2024),
24 most models focus on monaural audio and neglect spatialized soundscapes that enhance immersion.
25 The human binaural system uses interaural time differences (ITD) and interaural level differences
26 (ILD) to localize sound in three-dimensional space, and spatial audio must remain consistent with
27 visual cues. Any mismatch can disrupt immersion and weaken the sense of presence. For example,
28 hearing a cat's meow from the left immediately suggests its location, even if it is just off-screen.

29 Despite the growing importance of spatial audio in these immersive technologies, progress in machine
30 learning for spatial audio understanding is limited by the fundamental spatial data constraints. Most
31 existing audio datasets(Gemmeke et al., 2017; Chen et al., 2020; Agostinelli et al., 2023) focus on
32 monaural recordings, which discard vital spatial information, effectively "flattening" the soundscape
33 and preventing models from learning key physical phenomena such as room reverberation, echo
34 patterns, and sound propagation. Moreover, the scarcity of multimodal datasets that align spatial audio
35 with synchronized visual, position geometric, and semantic annotations hinders the development of
36 advanced auditory scene-analysis systems capable of human-like spatial perception.

37 To address these gaps, we present **MRSAudio**, a 500-hour large-scale multimodal spatial audio
38 dataset designed to support both spatial audio understanding and generation. It integrates high-fidelity
39 spatial recordings with synchronized video, 3D pose tracking, and rich semantic annotations, enabling
40 comprehensive modeling of real-world auditory scenes. As shown in Figure 1, the dataset comprises
41 four subsets, each targeting distinct tasks and scenarios. **MRSLife** (150 h) captures daily activities
42 such as board games, cooking, and office work, using egocentric video and FOA audio annotated with
43 sound events and speech transcripts. **MRSSpeech** (200 h) includes binaural conversations from 50
44 speakers across diverse indoor environments, paired with video, 3D source positions, and complete
45 scripts. **MRSSing** (75 h) features high-quality solo singing performances in Chinese, English,
46 German, and French by 20 vocalists, each aligned with time-stamped lyrics and corresponding
47 musical scores. **MRSMusic** (75 h) offers spatial recordings of 23 Traditional Chinese, Western and
48 Electronic instruments, with symbolic score annotations that support learning-based methods for
49 symbolic-to-audio generation and fine-grained localization. Together, these four subsets support a
50 broad spectrum of spatial audio research problems, including event detection, sound localization, and
51 binaural or ambisonic audio generation. By pairing spatial audio with synchronized exocentric and
52 egocentric video, geometric tracking, and detailed semantic labels, MRSAudio enables new research
53 directions in multimodal spatial understanding and cross-modal generation. Throughout this paper,
54 unless stated otherwise, we use the term spatial audio to refer to binaural audio.

- 55 • We introduce **MRSAudio**, a 500-hour, large-scale multimodal spatial audio dataset explicitly
56 designed to push the boundaries of spatial audio understanding and generative modeling.
- 57 • We assemble synchronized binaural and ambisonic recordings with exocentric and egocentric video,
58 geometric source positions, transcripts, scores, lyrics, and event labels, providing one of the most
59 richly annotated multimodal resources for spatial audio research.
- 60 • We organize the data into four complementary subsets (MRSLife, MRSSpeech, MRSSing, MRSMu-
61 sic), each carefully tailored to different real-world acoustic scenarios and equipped with rich,
62 scenario-specific annotations to facilitate downstream task development.
- 63 • We establish and release evaluation protocols and baseline implementations for five benchmark
64 tasks: audio spatialization generation, spatial text-to-speech, spatial singing voice synthesis, spa-
65 tial music generation, and localization and detection of sound events, in order to demonstrate
66 MRSAudio's versatility and to foster reproducible research.

67 The remainder of this paper is organized as follows. Section 2 reviews existing spatial audio datasets.
68 Section 3 describes the design, collection process, and key statistics of MRS Audio. Section 4 presents
69 extensive benchmark experiments using state-of-the-art methods on five core spatial audio tasks:
70 audio spatialization, spatial text-to-speech, spatial singing voice generation, spatial music synthesis,
71 and sound event localization and detection. Finally, Section 5 concludes the paper and discusses the
72 limitations and potential risks associated with MRS Audio.

73 **2 Related Work**

74 Deep learning has driven remarkable progress in both audio generation (Huang et al., 2023a;
 75 Dongchao Yang, 2023; Kreuk et al., 2022) and audio understanding(Chu et al., 2023; Huang et al.,
 76 2023b; Tang et al., 2024) tasks. However, the majority of these advances still rely heavily on monau-
 77 ral audio, which lacks the ability to represent or capture the rich spatial cues that naturally occur
 78 in real-world environments. The rapid adoption of VR/AR technologies has concurrently driven
 79 growing demand for immersive spatial audio experiences. Researchers have focused on several key
 80 technologies including: sound event localization and detection (Adavanne et al., 2019; Wang et al.,
 81 2022), mono-to-spatial audio conversion (Gao & Grauman, 2019; Pedro Morgado & Wang, 2018),
 82 and end-to-end spatial audio generation (Sun et al., 2024; Kim et al., 2025; Liu et al., 2025).

83 Despite recent progress, spatial audio generation and understanding remain constrained by the paucity
 84 of high-quality datasets. Due to the difficulty and expense of collecting and annotating high-quality
 85 spatial audio datasets, most open-source datasets still primarily consist of simulated or web-crawled
 86 content. Simulated datasets offer precise annotations but lack perceptual realism, while crawled
 87 datasets often provide real-world diversity but are missing critical labels, such as listener and source
 88 positions and content-level annotations, thus limiting their utility. For instance, spatial symbolic
 89 music generation demands accurate musical scores and corresponding positional metadata. To better
 90 understand the current landscape, we survey existing spatial audio datasets. These datasets differ in
 91 terms of audio format, including First-Order Ambisonics (FOA), multi-channel arrays, and binaural
 92 microphones, as well as in their collection methods (simulated, crawled, recorded) and annotation.

Table 1: Comparison of spatial audio datasets, where T denotes speech transcripts, P represents sound source positions, N indicates natural language prompts, and C stands for sound class tags

Dataset	Audio Format			Collect	Hours	Type	Visual	Label			
	FOA	Multi	Binaural					T	P	N	C
Spatial LibriSpeech	✓	✓	✗	Simulated	650	Speech	-	✓	✓	✗	✗
YT-Ambigen	✓	✗	✗	Crawled	142	ALL	Video	✗	✗	✗	✗
BEWO-1M	✗	✗	✓	Craw+Sim	2800	ALL	Image	✗	✗	✓	✗
FAIR-Play	✗	✗	✓	Recorded	5.2	Music	Video	✗	✗	✗	✗
STARSS23	✓	✓	✗	Recorded	7.5	ALL	Video	✗	✓	✗	✓
BinauralMusic	✗	✗	✓	Crawled	15.2	Music	Video	✗	✗	✗	✓
Sphere360	✓	✗	✗	Crawled	288	ALL	Video	✗	✗	✗	✗
MRSAudio (Ours)	✓	✗	✓	Recorded	500	ALL	Video	✓	✓	✓	✓

93 As shown in Table 1, existing datasets vary in their goals and modalities. For instance, Spatial
 94 LibriSpeech (Sarabia et al., 2023) simulates spatial speech using the LibriSpeech corpus and is mainly
 95 intended for binaural TTS applications. YT-Ambigen (Kim et al., 2025) and Sphere360 (Liu et al.,
 96 2025) are derived from web-crawled video datasets, but lack explicit spatial annotations, making them
 97 suitable primarily for video-to-spatial audio generation. BinauralMusic focuses on musical content,
 98 offering instrument class tags. BEWO-1M (Sun et al., 2024) combines crawled content and simulated
 99 binaural rendering, and provides image or GPT-generated prompts. STARSS23(Shimada et al.,
 100 2023) features real-world FOA and multi-channel audio alongside synchronized videos and includes
 101 sound class and position labels. In contrast to prior datasets, MRSAudio offers a comprehensive,
 102 large-scale, and real-world spatial audio corpus, featuring over 500 hours of recorded data in both
 103 FOA and binaural formats, covering all audio domains including general audio, speech, singing, and
 104 music. Uniquely, MRSAudio includes synchronized audio, video, and 3D positional geometry, along
 105 with fine-grained cross-modal annotations, such as: transcripts, word and phoneme boundaries and
 106 music scores. These features make MRSAudio an ideal resource for a broad range of spatial audio
 107 generation and understanding tasks, including spatial speech, music, and singing voice synthesis.

108 **3 Dataset Description**

109 In this section, we present **MRSAudio**, a freely available multimodal spatial audio corpus with
 110 synchronized video, positional data, and fine-grained annotations, released under the CC BY-NC-SA
 111 4.0 license. Figure 2 illustrates our data processing pipeline, with detailed descriptions in subsequent
 112 subsections. We then summarize key statistics that demonstrate MRSAudio’s scale and diversity.

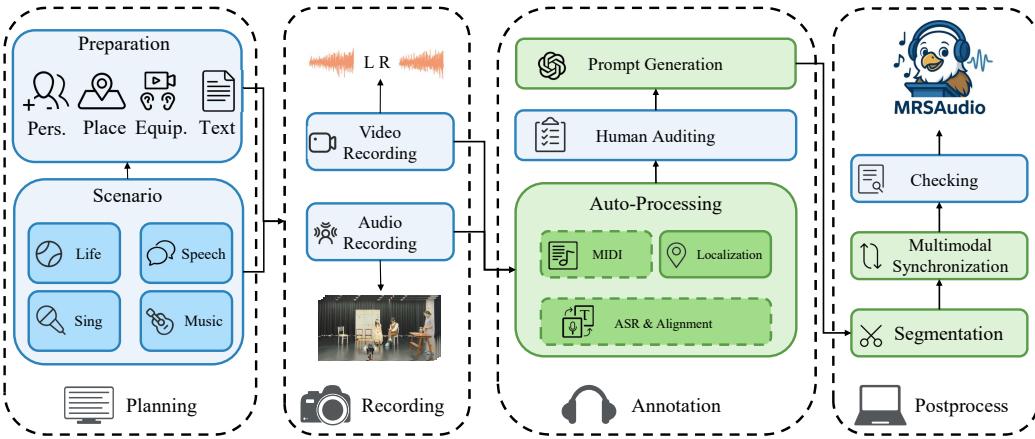


Figure 2: The pipeline of data collection and processing of MRSAudio. The blue boxes indicate steps requiring manual intervention, while the green boxes denote automated processing. In the “auto-processing” section, dashed modules apply to some scenarios, while solid modules apply to all.

113 3.1 Planning

114 To ensure that MRSAudio comprehensively covers scenarios from daily life, speech, singing, and
 115 music, we conduct systematic and modular planning before recording as follows.

116 **MRSLife:** This subset focuses on everyday conversations and environmental sound events. Based on
 117 the degree of human vocal interaction, MRSLife is further divided into two parts: MRSDialogue,
 118 which captures unscripted conversations that naturally include spontaneous action sounds (e.g.,
 119 footsteps, door movements), and MRSSound, which focuses on non-verbal sound events primarily
 120 caused by physical activities such as cooking, typing, or sports.

121 **MRSSpeech:** MRSSpeech targets clean, high-quality conversational recordings for TTS task. All
 122 spoken interactions are recorded in controlled indoor environments with minimal noise. We invite
 123 speakers to participate in content-driven conversations based on predefined scripts.

124 **MRSSing:** MRSSing captures solo vocal performances for singing voice synthesis tasks. It includes
 125 professional solo vocal recordings in four languages: Mandarin, English, German, and French,
 126 performed by singers covering the full vocal range including soprano, alto, tenor, and bass.

127 **MRSMusic:** MRSMusic captures immersive instrumental performances suitable for spatial music
 128 generation and analysis. We record solo performances from 45 professional musicians across 23
 129 instruments, including Traditional Chinese, Western and Electronic instruments. Each performance is
 130 paired with its corresponding musical score to support score-based music generation.

131 Details regarding personnel recruitment, venue selection, equipment configuration, and material
 132 preparation for each subset are provided in Appendix A.1.

133 3.2 Recording

134 Based on the predefined planning, we conduct parallel data collection across all modules. To ensure
 135 participant anonymity, masks are worn during recording when necessary. All participants sign an
 136 open-source data release agreement, allowing the dataset to be freely distributed for academic research
 137 purposes. The recording details are summarized as follows:

138 **MRSLife:** In MRSDialogue, Audio is recorded using a professional binaural recording head and
 139 high-resolution sound cards, while synchronized video is captured using industry-standard cameras.
 140 In MRSSound, in addition to binaural audio and exocentric video, we also captured FOA (Zoom
 141 H3-VR) and egocentric video (Gopro camera). To ensure the effectiveness of the egocentric video,
 142 participants are asked to remain within the frontal field of view of the binaural recording head.

143 **MRSSpeech:** To introduce spatial variability, we select four recording rooms that differ in size and
144 acoustic material. Each speaker reads inflected emotions from the scripts while walking through
145 the space, producing dynamic spatial cues. Recordings include spoken passages, binaural audio and
146 exocentric video, as well as a clean mono audio by a lavalier microphone placed near the speaker.

147 **MRSSing:** The professional singers perform according to musical scores. To introduce variation in
148 source–listener geometry, we adjust the position of the head-mounted binaural microphone relative to
149 the singer. In addition to the binaural audio and synchronized video, each session includes a clean
150 vocal track recorded with a studio-grade condenser microphone.

151 **MRSMusic:** We record 23 Traditional Chinese, Western and Electronic instruments, performed by
152 45 professional musicians. To capture rich spatial detail, we vary microphone placement around the
153 instrument. Recordings include binaural audio, monaural audio, exocentric video, and synchronized
154 video that records playing gestures. Full recording details are provided in Appendix A.2.

155 3.3 Annotation

156 To maximize MRSAudio’s utility across a wide range of tasks, we begin with event-level annotations
157 for all vocal and acoustic content in MRSLife, MRSSpeech, MRSSing, and MRSMusic. However,
158 coarse annotations alone are insufficient for fine-grained tasks such as singing voice modeling and
159 music generation from scores. To bridge this gap, we design a comprehensive annotation pipeline.
160 Full implementation details and detailed annotation guidelines are provided in Appendix A.3.

161 **MRSLife:** For MRSDialogue, we apply WhisperX (Bain et al., 2023) for automatic speech recogni-
162 tion and speaker diarization to generate initial transcripts and speaker turns. Human annotators
163 correct recognition errors and speaker attribution mismatches. The audio is then segmented into
164 utterances and the transcripts are converted into phoneme sequences (using pypinyin for Mandarin). A
165 two-stage alignment process follows: we first apply the Montreal Forced Aligner (MFA) (McAuliffe
166 et al., 2017) for coarse word/phoneme mapping, then manually refine boundaries in Praat (Boersma,
167 2001). For MRSSound, we annotate sound event categories and corresponding time intervals.

168 **MRSSpeech:** Given the availability of full scripts, we adapt WhisperX for long-form word-to-audio
169 alignment of up to 30 minutes. Each script line is automatically matched to its corresponding audio
170 segment (see Appendix A.3 for more details). Annotators then review these alignments, correcting
171 any omissions or insertions caused by actors’ deviations from the script. Finally, phoneme sequences
172 are extracted and aligned with the audio using the same procedure as in MRSDialogue.

173 **MRSSing:** We use voice activity detection (VAD) to segment recordings into singing regions, then
174 align pre-existing lyrics using LyricFA’s ASR-based dynamic programming algorithm³. Phoneme
175 generation is language-dependent: pypinyin for Mandarin, ARPA for English, and MFA’s built-in
176 phoneme sets for French and German. Alignment is conducted via MFA, followed by manual
177 refinement. Melody and rhythm of singing are transcribed into MIDI format using ROSVOT (Li
178 et al., 2024). Annotators then label each excerpt with high-level style descriptors, such as emotional
179 tone (happy, sad), tempo (slow, moderate, fast), and pitch range (low, medium, high).

180 **MRSMusic:** We use Audio Slicer⁴ to segment the music recordings and generate initial symbolic
181 annotations with basic-pitch (Bittner et al., 2022), and then employ professional musicians to verify
182 and adjust note onsets, offsets, and dynamics to match the performance accurately.

183 **Source Localization:** For static sources, we manually record 3D positions relative to the capture
184 space. For dynamic scenes, we use the Ultra-Wideband (UWB) system(Aiello & Rogerson, 2003)
185 to track the positions of sound sources in real time. Based on the recorded position trajectories, we
186 generate natural language motion descriptions using GPT-4o (Achiam et al., 2023).

187 3.4 Post-Processing

188 The post-processing pipeline comprises three key steps to refine raw annotated data. First, segmen-
189 tation splits continuous recordings into task-oriented clips: utterances for speech and singing are
190 extracted via alignment timestamps, while MRSSound audio is uniformly divided into 10-second
191 segments. Next, multimodal synchronization aligns each clip with auxiliary modalities (text, video,

³<https://github.com/wolfgitpr/LyricFA>

⁴<https://github.com/flutydeer/audio-slicer>

192 position metadata) with temporal anchors. Static sound sources are annotated with manually measured 3D coordinates, whereas dynamic sources leverage interpolated UWB tracking trajectories.
 193 For scenes where participants' faces are visible, anonymization is performed by adding half-face
 194 masks. Finally, quality assurance involves domain experts auditing 15% of clips across all modules
 195 to verify temporal alignment precision, cross-modal content consistency, and annotation accuracy.
 196 Full auditing protocols are documented in Appendix A.4, ensuring reproducibility of this process.
 197

198 3.5 Statistics

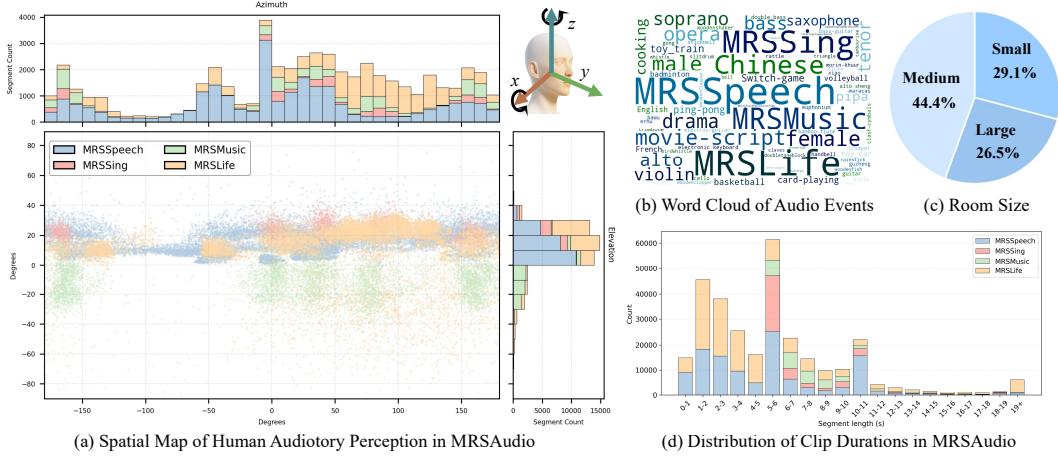


Figure 3: Statistical overview of MRSAudio. (a) Spatial distribution of sound sources relative to the listener. Red, green, and blue arrows denote the positive x, y, and z axes; azimuth is measured around the z-axis from the x-axis, and elevation is relative to the xy-plane. (b) Word cloud. (c) Proportions of recording spaces by room size. (d) Distribution of audio segment durations.

199 To illustrate spatial diversity, Figure 3(a) shows the 3D distribution of sound source positions
 200 relative to the listener. The heatmap reveals near-uniform azimuthal coverage, with greater density
 201 in the frontal hemisphere due to the prevalence of egocentric video recordings. Elevation angles
 202 are concentrated between -40° and 40° , aligning with typical human sound perception patterns.
 203 Figure 3(b) summarizes the annotations using a keyword cloud that captures the range of real-world
 204 activities recorded. Figure 3(c) presents the proportions of recording spaces by room size: medium-
 205 sized rooms are the most common, accounting for approximately 40% of sessions, while small and
 206 large rooms each represent around 30%. Recording duration is evenly distributed across room types,
 207 ensuring diverse acoustic coverage. Figure 3(d) shows the distribution of segment durations after
 208 automatic and manual segmentation. Most audio clips are shorter than 10 seconds, which is suitable
 209 for modeling short-duration events and supports efficient downstream training and inference. Overall,
 210 MRSAudio offers comprehensive coverage across spatial positions, acoustic environments, scene
 211 types, and temporal structures, making it well-suited for a wide range of spatial audio generation and
 212 understanding tasks. A more detailed breakdown of per-scenario statistics is provided in Appendix A.

213 4 Benchmarks

214 To demonstrate the quality and utility of MRSAudio in real-world scenarios, we evaluate it on five
 215 representative spatial audio tasks: (i) audio spatialization, (ii) spatial text-to-speech, (iii) spatial
 216 singing voice synthesis, (iv) spatial music generation, and (v) sound event localization and detection
 217 (SELD). These tasks cover both generation and understanding, and are critical for applications such as
 218 AR/VR, spatial media production, and perceptual scene analysis. All experiments are conducted using
 219 state-of-the-art methods on a server equipped with eight NVIDIA RTX 4090 GPUs. For different
 220 tasks, we employ distinct objective metrics. For generation tasks, we compute cosine similarity scores
 221 for direction (ANG Cos) and distance (DIS Cos) to quantify spatial alignment quality. Additionally,
 222 we utilize subjective MOS-Q (Mean Opinion Score for Quality) to evaluate the quality of generated
 223 audio and MOS-P (Mean Opinion Score for Position) to assess spatial perception. For implementation
 224 details of training and evaluation metrics, please refer to Appendix B.

225 **4.1 Audio Spatialization**

226 Audio spatialization aims to synthesize spatially immersive soundscapes from monaural inputs and
 227 source positional information. To evaluate MRSAudio on this task, we adopt BinauralGrad (Leng
 228 et al., 2022), a diffusion-based generation model that predicts binaural waveforms conditioned on
 229 monaural audio and source coordinates. Since the downstream generation tasks can be formulated as
 230 predicting interaural (binaural) representations from mono audio, we employ a single-stage training
 231 scheme for all experiments. For comparison, we include a traditional signal processing baseline
 232 (DSP), which renders binaural audio using virtual source positions simulated via room impulse
 233 responses (RIRs) and head-related transfer functions (HRTFs). We adopt the following objective
 234 metrics to evaluate audio quality: (1) W-L2: waveform L2 distance, (2) A-L2: amplitude envelope L2
 235 distance, (3) P-L2: phase difference L2, (4) STFT: multi-resolution STFT loss, (5) PESQ: perceptual
 236 speech quality score. Results for the ground truth and DSP baseline are obtained by averaging across
 237 all test sets from the four MRSAudio subsets. Further details are provided in Appendix B.3.

Table 2: Audio Spatialization Performance. For MRSLife, we only use the MRSSound subset.

Method	W-L2 $\times 10^{-3} \downarrow$	A-L2 \downarrow	P-L2 \downarrow	PESQ \uparrow	STFT \downarrow	MOS-Q \uparrow	MOS-P \uparrow
Ground Truth	-	-	-	-	-	4.69 \pm 0.08	4.56 \pm 0.10
DSP	1.691	0.048	1.562	2.830	1.246	3.89 \pm 0.09	3.75 \pm 0.11
MRSLife	0.076	0.025	0.898	-	2.243	3.91 \pm 0.07	3.87 \pm 0.10
MRSSpeech	0.460	0.061	0.807	1.929	2.352	3.88 \pm 0.08	3.84 \pm 0.08
MRSSing	0.647	0.093	1.004	1.723	2.539	3.84 \pm 0.09	3.91 \pm 0.07
MRSMusic	0.705	0.063	0.835	-	1.724	3.87 \pm 0.07	3.93 \pm 0.09

238 As shown in Table 2, BinauralGrad achieves strong performance across all MRSAudio subsets, sur-
 239 passing the classical DSP baseline on most objective and subjective metrics. The highest performance
 240 is observed on MRSLife, likely due to the relatively limited variation in sound sources within these
 241 scenes. These results demonstrate that MRSAudio’s rich spatial annotations and diverse acoustic
 242 environments provide a solid foundation for training and evaluating spatial audio generation models.

243 **4.2 Spatial Text to Speech**

244 Spatial TTS aims to produce high-quality speech enriched with spatial cues, thereby enhancing
 245 immersion and realism in AR/VR. Although recent advances in TTS have led to impressive im-
 246 provements in speech quality, progress in spatialized speech generation remains limited due to the
 247 scarcity of spatially annotated recordings with rich, high-quality labels. To evaluate the effectiveness
 248 of MRSAudio for this task, we train an end-to-end model following ISDrama (Zhang et al., 2025) to
 249 directly generate spatial speech from text and position data. Additionally, we compare with cascaded
 250 pipelines that combine monaural TTS models (CosyVoice (Du et al., 2024) and F5-TTS (Chen et al.,
 251 2024)) with the audio spatialization module. This allows us to assess both speech generation quality
 252 and the spatial fidelity enabled by our dataset. For content evaluation, we report Character Error Rate
 253 (CER) and Speaker Similarity (SIM). Further details are provided in Appendix B.4.

Table 3: Spatial TTS Performance on MRSSpeech. “SP” denotes the Audio Spatialization.

Method	Objective				Subjective	
	CER \downarrow	SIM \uparrow	ANG Cos \uparrow	DIS Cos \uparrow	MOS-Q \uparrow	MOS-P \uparrow
Ground Truth	2.54%	-	-	-	4.39 \pm 0.08	4.16 \pm 0.10
Mono + SP	2.56%	0.98	0.44	0.68	3.88 \pm 0.08	3.84 \pm 0.08
CosyVoice + SP	3.89%	0.96	0.41	0.63	3.75 \pm 0.12	3.72 \pm 0.09
F5-TTS + SP	3.15%	0.97	0.40	0.62	3.69 \pm 0.13	3.67 \pm 0.14
ISDrama (speech)	3.35%	0.96	0.48	0.65	3.85 \pm 0.09	3.82 \pm 0.11

254 As shown in Table 3, the Mono+SP method achieves strong performance across most metrics. The
 255 CER remains low and comparable to the ground truth, indicating preserved linguistic accuracy after
 256 spatialization. A high SIM score reflects stable timbre learning. ANG Cos and DIS Cos show
 257 good spatial alignment with the ground truth, and subjective MOS scores confirm that the generated

258 speech is both natural and spatially coherent. These results demonstrate that MRSSpeech provides
 259 high-quality, spatially annotated training data that enables effective and realistic spatial TTS.

260 4.3 Spatial Singing Voice Synthesis

261 Spatial SVS aims to produce expressive, high-quality singing voices enriched with accurate spatial
 262 cues, thereby enhancing listener immersion. While traditional SVS has advanced considerably, spatial
 263 SVS remains underexplored due to the lack of high-quality, spatially annotated datasets. To evaluate
 264 MRSpeech for this task, we use the MRSSing subset to train and benchmark models. We adopt
 265 the ISDrama architecture to perform end-to-end spatial singing synthesis, incorporating note-level
 266 pitch control to enhance prosody accuracy. For comparison, we use the open-source SVS models
 267 Rmssinger (He et al., 2023). The outputs are then spatialized with BinauralGrad. We employ objective
 268 metrics, including Mel-Cepstral Distortion (MCD) and F0 Frame Error (FFE), to evaluate spectral
 and pitch similarity between predicted and ground-truth. Details are in Appendix B.5.

Table 4: Spatial SVS Performance on MRSSing. “SP” denotes the Audio Spatialization.

Method	Objective				Subjective	
	MCD ↓	FFE ↓	ANG Cos ↑	DIS Cos ↑	MOS-Q ↑	MOS-P ↑
GT (Ground Truth)	-	-	-	-	4.45 ± 0.10	4.30 ± 0.12
Mono+SP	3.19	0.17	0.51	0.71	3.84 ± 0.09	3.91 ± 0.07
Rmssinger+SP	3.85	0.23	0.45	0.65	3.65 ± 0.08	3.81 ± 0.11
ISDrama(sing)	3.71	0.21	0.47	0.70	3.86 ± 0.13	3.88 ± 0.09

269
 270 As shown in Table 4, the Mono + SP approach trained on MRSSing with pitch control achieves the
 271 best performance across most metrics. Its low MCD indicates strong spectral fidelity, and high ANG
 272 Cos and DIS Cos scores demonstrate effective spatial alignment. Subjective MOS results confirm
 273 that the generated singing voices are natural, high-quality, and spatially coherent. These findings
 274 validate MRSSing as an effective resource for spatial singing voice generation.

275 4.4 Spatial Music Generation

276 This task aims to synthesize spatially immersive music conditioned on symbolic scores. While datasets
 277 like FAIR-Play offer high-quality instrument recordings, they lack aligned sheet music, limiting
 278 controllable generation. In contrast, our MRSMusic subset includes 23 instruments with aligned
 279 scores, enabling fine-grained, score-based spatial music synthesis. We benchmark three systems on
 280 MRSMusic. First, we apply BinauralGrad (Leng et al., 2022) to spatialize mono recordings. Second,
 281 we utilize Make-An-Audio 2 (Copet et al., 2023) to generate mono music from MIDI scores, which
 282 is subsequently spatialized. Third, we adapt ISDrama to accept both score embeddings and spatial
 283 poses, enabling end-to-end spatial symbolic music generation. For objective evaluation, we compute
 284 Fréchet Audio Distance (FAD) and FFE to evaluate the results. Details are in Appendix B.6.

Table 5: Spatial MG Performance on MRSMusic. “SP” denotes the Audio Spatialization.

Method	Objective				Subjective	
	FAD ↓	FFE ↓	ANG Cos ↑	DIS Cos ↑	MOS-Q ↑	MOS-P ↑
GT (Ground Truth)	-	-	-	-	4.49 ± 0.09	4.34 ± 0.11
Mono+SP	2.88	0.14	0.48	0.74	3.87 ± 0.07	3.93 ± 0.09
Make-An-Audio 2+SP	4.39	0.49	0.49	0.41	3.74 ± 0.10	3.73 ± 0.13
ISDrama(music)	2.45	0.21	0.53	0.68	3.89 ± 0.12	3.88 ± 0.10

285 As shown in Table 5, the Mono + SP pipeline achieves strong performance, benefiting from access to
 286 ground truth mono audio. The Make-An-Audio 2 + SP exhibits limitations in FAD and pitch accuracy.
 287 This suggests that general-purpose audio generation models may struggle to fully capture structured
 288 musical information from symbolic prompts. The ISDrama (Music) model generates music directly
 289 from the symbolic inputs and spatial cues, achieving better coherence and spatial alignment than the
 290 Make-An-Audio 2 + SP. These results highlight MRSMusic’s effectiveness in supporting spatially
 291 controllable music generation across diverse instruments and spatial conditions.

292 **4.5 Sound Event Localization and Detection**

293 This task evaluates the ability to detect and localize sound events using MRSAudio’s spatial annotations.
 294 We follow STARSS23 (Shimada et al., 2023) using both audio-only and audio-visual
 295 variants on the MRSSound subset. In the audio-only condition, models receive either FOA or binaural
 296 waveforms as input and predict sound event classes along with 3D source coordinates. For the
 297 audio-visual condition, we extract bounding boxes of visible persons to serve as coarse visual priors,
 298 which are fused with the audio representations. We also explore architectural variations by replacing
 299 the original convolutional backbone with a Transformer encoder, allowing us to assess the impact
 300 of temporal modeling capacity on spatial prediction. We evaluate model performance using four
 301 standard joint detection and localization metrics, including location-aware detection (F_{20° , ER_{20°)
 302 and class-aware localization (LE_{CD} , LR_{CD}). Details are in Appendix B.7.

Table 6: Sound Event Localization and Detection on MRSSound.

Model	Audio Type	Visual	$ER_{20^\circ} \downarrow$	$F_{20^\circ} \uparrow$	$LE_{CD} \downarrow$	$LR_{CD} \uparrow$
ConvNet	FOA	✓	1.17 ± 0.02	13.00 ± 2.72	42.44 ± 8.91	82.76 ± 5.19
ConvNet	FOA	✗	1.12 ± 0.02	9.33 ± 0.36	46.47 ± 4.46	85.35 ± 3.80
ConvNet	Binaural	✓	1.11 ± 0.03	5.90 ± 3.86	41.95 ± 9.85	45.81 ± 11.29
ConvNet	Binaural	✗	1.11 ± 0.02	6.00 ± 0.34	45.17 ± 7.92	70.14 ± 10.38
Transformer	FOA	✓	1.01 ± 0.01	7.52 ± 0.42	35.69 ± 7.47	46.78 ± 7.62
Transformer	FOA	✗	0.99 ± 0.05	7.95 ± 0.15	48.76 ± 2.66	87.32 ± 2.14
Transformer	Binaural	✓	1.01 ± 0.02	8.47 ± 0.65	36.18 ± 7.29	41.33 ± 12.73
Transformer	Binaural	✗	1.11 ± 0.04	7.37 ± 0.97	46.74 ± 7.91	43.87 ± 10.74

303 As shown in Table 6, model performance varies with architecture, input modality, and audio format.
 304 Transformer-based models generally outperform ConvNet baselines, particularly in reducing localiza-
 305 tion error. FOA input consistently yields better results than binaural audio, benefiting from richer
 306 spatial representation. For example, the Transformer with FOA and no visual input achieves the lowest
 307 error rate ($ER_{20^\circ} 0.99$) and highest localization recall ($LR_{CD} 87.32$). The addition of visual features
 308 improves performance in some ConvNet settings (e.g., F_{20° rises from 9.33 to 13.00 with FOA), but
 309 offers limited or inconsistent gains in Transformer models, possibly due to modality mismatch or
 310 redundancy. These results highlight the value of MRSAudio’s spatial annotations and multimodal
 311 streams in supporting flexible evaluation under both unimodal and multimodal configurations.

312 **5 Conclusion and Discussion**

313 We introduce MRSAudio, a large-scale, multimodal spatial audio corpus designed to support a
 314 wide range of generation and understanding tasks. MRSAudio comprises four complementary
 315 modules, MRSLife, MRSSpeech , MRSSing, and MRSMusic, each captured with binaural/FOA
 316 audio, synchronized video, precise 3D pose metadata, and richly detailed annotations (event labels,
 317 transcripts, phoneme boundaries, lyrics, musical scores, and motion prompts). Through extensive
 318 benchmarks on audio spatialization, binaural speech, singing, and music generation , as well as sound
 319 event localization and detection, we demonstrate that MRSAudio’s scale, diversity, and annotation
 320 depth enable state-of-the-art performance and unlock new avenues for spatial audio research.

321 **Limitations and Future Directions:** While MRSAudio offers broad multimodal coverage, two limi-
 322 tations remain. First, although synchronized video is provided for all recordings, current benchmarks
 323 only explore visual input in a limited subset of tasks. Second, while the dataset includes both binaural
 324 and First-Order Ambisonic (FOA) formats, the FOA subset is smaller in scale, and most tasks focus
 325 on binaural audio, limiting spatial modeling diversity. Future work will expand the role of visual
 326 modalities in tasks such as sound localization and scene understanding. We also plan to increase FOA
 327 recordings to balance data distribution and support broader spatial audio research, and develop more
 328 FOA-specific benchmarks to better utilize ambisonic spatial cues.

329 **Negative Societal Impact:** As with any large-scale audiovisual dataset, MRSAudio carries potential
 330 risks if misused. It could be exploited to generate highly realistic yet synthetic spatial audio for
 331 deepfakes or disinformation in AR and VR applications. To mitigate such risks, MRSAudio is
 332 released under a noncommercial license with clear usage guidelines. We encourage responsible use in
 333 accordance with ethical standards, including consent management, data governance, and transparency.

334 **References**

- 335 Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt,
336 J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 337 Adavanne, S., Politis, A., and Virtanen, T. A multi-room reverberant dataset for sound event
338 localization and detection. *Proc. DCASE2019*, 2019.
- 339 Agostinelli, A., Denk, T. I., Borsos, Z., Engel, J., Verzetti, M., Caillon, A., Huang, Q., Jansen,
340 A., Roberts, A., Tagliasacchi, M., et al. Musiclm: Generating music from text. *arXiv preprint
341 arXiv:2301.11325*, 2023.
- 342 Aiello, G. R. and Rogerson, G. D. Ultra-wideband wireless systems. *IEEE microwave magazine*, 4
343 (2):36–47, 2003.
- 344 Baevski, A., Zhou, Y., Mohamed, A., and Auli, M. wav2vec 2.0: A framework for self-supervised
345 learning of speech representations. *Advances in neural information processing systems*, 33:
346 12449–12460, 2020.
- 347 Bain, M., Huh, J., Han, T., and Zisserman, A. Whisperx: Time-accurate speech transcription of
348 long-form audio. *arXiv preprint arXiv:2303.00747*, 2023.
- 349 Bittner, R. M., Bosch, J. J., Rubinstein, D., Meseguer-Brocal, G., and Ewert, S. A lightweight
350 instrument-agnostic model for polyphonic note transcription and multipitch estimation. In *Proceed-
351 ings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*,
352 Singapore, 2022.
- 353 Boersma, P. Praat, a system for doing phonetics by computer. *Glot. Int.*, 5(9):341–345, 2001.
- 354 Chen, H., Xie, W., Vedaldi, A., and Zisserman, A. Vggsound: A large-scale audio-visual dataset. In
355 *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing
(ICASSP)*, pp. 721–725. IEEE, 2020.
- 357 Chen, Y., Niu, Z., Ma, Z., Deng, K., Wang, C., Zhao, J., Yu, K., and Chen, X. F5-tts: A fairytaler that
358 fakes fluent and faithful speech with flow matching. *arXiv preprint arXiv:2410.06885*, 2024.
- 359 Chu, Y., Xu, J., Zhou, X., Yang, Q., Zhang, S., Yan, Z., Zhou, C., and Zhou, J. Qwen-audio:
360 Advancing universal audio understanding via unified large-scale audio-language models. *arXiv
361 preprint arXiv:2311.07919*, 2023.
- 362 Copet, J., Kreuk, F., Gat, I., Remez, T., Kant, D., Synnaeve, G., Adi, Y., and Défossez, A. Simple
363 and controllable music generation. *arXiv preprint arXiv:2306.05284*, 2023.
- 364 Dongchao Yang, Jinchuan Tian, X. T. R. H. S. L. X. C. J. S. S. Z. J. B. X. W. Z. Z. H. M. Uniaudio:
365 An audio foundation model toward universal audio generation. *arXiv preprint arXiv:2310.00704*,
366 2023.
- 367 Du, Z., Chen, Q., Zhang, S., Hu, K., Lu, H., Yang, Y., Hu, H., Zheng, S., Gu, Y., Ma, Z., et al.
368 Cosyvoice: A scalable multilingual zero-shot text-to-speech synthesizer based on supervised
369 semantic tokens. *arXiv preprint arXiv:2407.05407*, 2024.
- 370 Gao, R. and Grauman, K. 2.5d visual sound. *arXiv preprint arXiv:1812.04204*, 2019.
- 371 Gao, Z., Li, Z., Wang, J., Luo, H., Shi, X., Chen, M., Li, Y., Zuo, L., Du, Z., Xiao, Z., et al. Funasr:
372 A fundamental end-to-end speech recognition toolkit. *arXiv preprint arXiv:2305.11013*, 2023.
- 373 Gemmeke, J. F., Ellis, D. P. W., Freedman, D., Jansen, A., Lawrence, W., Moore, R. C., Plakal, M.,
374 and Ritter, M. Audio set: An ontology and human-labeled dataset for audio events. In *Proc. IEEE
375 ICASSP 2017*, New Orleans, LA, 2017.
- 376 He, J., Liu, J., Ye, Z., Huang, R., Cui, C., Liu, H., and Zhao, Z. Rmssinger: Realistic-music-score
377 based singing voice synthesis. *arXiv preprint arXiv:2305.10686*, 2023.
- 378 Huang, J., Ren, Y., Huang, R., Yang, D., Ye, Z., Zhang, C., Liu, J., Yin, X., Ma, Z., and Zhao, Z.
379 Make-an-audio 2: Temporal-enhanced text-to-audio generation, 2023a.

- 380 Huang, R., Li, M., Yang, D., Shi, J., Chang, X., et al. AudioGPT: Understanding and generating
 381 speech, music, sound, and talking head. *arXiv preprint arXiv:2304.12995*, 2023b.
- 382 Kilgour, K., Zuluaga, M., Roblek, D., and Sharifi, M. Fr\echet audio distance: A metric for
 383 evaluating music enhancement algorithms. *arXiv preprint arXiv:1812.08466*, 2018.
- 384 Kim, J., Yun, H., and Kim, G. Visage: Video-to-spatial audio generation. In *ICLR*, 2025.
- 385 Kreuk, F., Synnaeve, G., Polyak, A., Singer, U., Défossez, A., Copet, J., Parikh, D., Taigman, Y., and
 386 Adi, Y. Audiogen: Textually guided audio generation. *arXiv preprint arXiv:2209.15352*, 2022.
- 387 Leng, Y., Chen, Z., Guo, J., Liu, H., Chen, J., Tan, X., Mandic, D., He, L., Li, X., Qin, T., et al.
 388 Binauralgrad: A two-stage conditional diffusion probabilistic model for binaural audio synthesis.
 389 *Advances in Neural Information Processing Systems*, 35:23689–23700, 2022.
- 390 Li, R., Zhang, Y., Wang, Y., Hong, Z., Huang, R., and Zhao, Z. Robust singing voice transcription
 391 serves synthesis, 2024.
- 392 Liu, H., Luo, T., Jiang, Q., Luo, K., Sun, P., Wan, J., Huang, R., Chen, Q., Wang, W., Li, X., Zhang,
 393 S., Yan, Z., Zhao, Z., and Xue, W. Omniaudio: Generating spatial audio from 360-degree video,
 394 2025. URL <https://arxiv.org/abs/2504.14906>.
- 395 Liu, Z., Rosen, E., et al. Autoslicer: Scalable automated data slicing for ml model analysis. *arXiv
 396 preprint arXiv:2212.09032*, 2022.
- 397 McAuliffe, M., Socolof, M., Mihuc, S., Wagner, M., and Sonderegger, M. Montreal forced aligner:
 398 Trainable text-speech alignment using kaldi. In *Interspeech*, volume 2017, pp. 498–502, 2017.
- Pedro Morgado, Nuno Vasconcelos, T. L. and Wang, O. Self-supervised generation of spatial audio
 for 360deg video. In *Neural Information Processing Systems (NIPS)*, 2018.
- 399 Sarabia, M., Menyaylenko, E., Toso, A., Seto, S., Aldeneh, Z., Pirhosseinloo, S., Zappella, L., Theobald,
 400 B.-J., Apostoloff, N., and Sheaffer, J. Spatial librispeech: An augmented dataset for spatial audio
 401 learning. *arXiv preprint arXiv:2308.09514*, 2023.
- 402 Shimada, K., Politis, A., Sudarsanam, P., Krause, D. A., Uchida, K., Adavanne, S., Hakala, A., Koyama,
 403 Y., Takahashi, N., Takahashi, S., et al. Starss23: An audio-visual dataset of spatial recordings of real
 404 scenes with spatiotemporal annotations of sound events. *Advances in neural information processing
 405 systems*, 36:72931–72957, 2023.
- 406 Sun, P., Cheng, S., Li, X., Ye, Z., Liu, H., Zhang, H., Xue, W., and Guo, Y. Both ears wide open:
 407 Towards language-driven spatial audio generation. *arXiv preprint arXiv:2410.10676*, 2024.
- 408 Tang, C., Yu, W., Sun, G., Chen, X., Tan, T., et al. SALMONN: Towards generic hearing abilities for
 409 large language models. *Proc. ICLR*, 2024.
- 410 Wang, Q., Chai, L., Wu, H., Nian, Z., Niu, S., Zheng, S., Wang, Y., Sun, L., Fang, Y., Pan, J., et al. The
 411 nerc-slip system for sound event localization and detection of dcase2022 challenge. *DCASE2022
 412 Challenge, Tech. Rep.*, 2022.
- 413 Wang, Y., Guo, W., Huang, R., Huang, J., Wang, Z., You, F., Li, R., and Zhao, Z. Frieren: Efficient
 414 video-to-audio generation with rectified flow matching. *arXiv e-prints*, pp. arXiv–2406, 2024.
- 415 Wei, H., Cao, X., Dan, T., and Chen, Y. Rmvppe: A robust model for vocal pitch estimation in polyphonic
 416 music. *arXiv preprint arXiv:2306.15412*, 2023.
- 417 Zhang, Y., Guo, W., Pan, C., Zhu, Z., Jin, T., and Zhao, Z. Isdrama: Immersive spatial drama generation
 418 through multimodal prompting. *arXiv preprint arXiv:2504.20630*, 2025.
- 419 Zheng, Z., Peng, P., Ma, Z., Chen, X., Choi, E., and Harwath, D. Bat: Learning to reason about spatial
 420 sounds with large language models. *arXiv preprint arXiv:2402.01591*, 2024.

421 **A Details of Dataset**

422 **A.1 Details of Planning**

423 During the planning phase, we systematically analyze real-world application scenarios for spatial
424 audio and divide them into four major categories. MRSLife focuses on daily environmental sound
425 events, MRSSpeech targets high-quality conversational data, MRSSing captures solo vocal perfor-
426 mances, and MRSMusic records immersive instrumental music. These four scenarios collectively
427 cover a broad range of everyday acoustic contexts and are designed to support a wide spectrum of
428 downstream tasks.

429 Based on predefined scenario requirements, we recruit professionals from relevant fields to participate
430 in the recording process. To ensure spatial diversity, we rent various venues tailored to each scenario
431 type. Following the FAIR-Play protocol, we use a 3Dio Free Space XLR binaural microphone to
432 capture spatial audio, GoPro HERO cameras and mobile phones to record exocentric and egocentric
433 videos, and UWB-based tracking systems to capture motion trajectories. Once personnel, locations,
434 and equipment are secured, we prepare task-specific materials for each subset. For MRSLife, we
435 further divide the content based on the proportion of speech involved, resulting in two subcategories:
436 **MRSDialogue** (e.g., group games, board games) and **MRSSound** (e.g., kungfu, kitchens, offices,
437 sports). We predefine common sound events for each, such as clattering dishes, keyboard typing,
438 and table tennis. For MRSSpeech, we compile a large corpus of scripts from movies, TV shows,
439 and crosstalk performances and automatically extract dialogue passages for speaker delivery. For
440 MRSSing, we design lyrics in four languages (Chinese, English, German, and French) and recruit
441 singers across a range of vocal types to maximize diversity. For MRSMusic, we collect solo
442 performances across 23 traditional and modern instruments, covering a wide array of timbres and
443 playing techniques. Specific sound categories for each subset are summarized in Table 7.

Table 7: Examples of predefined content types for each MRSAudio subset.

Subset	Samples of Categories or Keywords
MRSLife	MRSDialogue: board games, card games, collaborative tasks MRSSound: kungfu, clattering dishes, cutting vegetables, typing, running, table tennis
MRSSpeech	Movie scripts, crosstalk, scripted TV dialogue, multi-speaker conversations
MRSSing	Chinese, English, German, French; soprano, alto, tenor, bass
MRSMusic	violin, electronic keyboard, cello, viola, double bass, trumpet, trombone, euphonium, erhu, pipa, xiao, bawu, trumpet, trombone, and others (23 instruments in total)

444 **A.2 Details of Recording**

445 We recruit a large number of participants with professional backgrounds in singing, music, and
446 language to contribute to the recording process. To protect their identity, participants are asked
447 to wear masks in appropriate scenarios. Prior to participation, all individuals sign consent forms
448 agreeing to the open-source release of their audio and video under the CC BY-NC-SA 4.0 license.

449 All audio is recorded in WAV format at a sampling rate of 48 kHz. Video is recorded at a minimum
450 resolution of 1080p and 24 frames per second, and is later standardized to this format during
451 post-processing.

452 **MRSLife:** We recruit 62 participants to perform daily activities including board games, cooking,
453 exercise, and office work. In MRSDialogue, each participant is compensated \$30 per recorded hour.
454 Binaural audio is captured using head-mounted microphones, and third-person video is recorded. In
455 MRSSound scenes such as kung fu or kitchen demonstrations, performers are paid \$50 per hour. Both
456 binaural and FOA (First-Order Ambisonics) recordings are collected, along with first-person and
457 third-person video. Participants receive brief scene descriptions and are asked to act naturally while
458 maintaining a single active sound source when possible. The total duration for MRSLife recordings
459 reaches approximately 150 hours.

460 **MRSSpeech:** We employ 50 expressive speakers to read from scripted texts, each paid \$30 per
461 hour of recorded audio. The total recorded duration reaches 200 hours, with compensation totaling

462 \$40,000. During recording sessions, speakers alternate between standing and walking around the
463 recording area to introduce spatial diversity while maintaining speech clarity. Binaural audio is
464 captured using head-mounted microphones, and clean monaural speech is recorded using lavalier
465 microphones. All sessions are filmed from a third-person perspective.

466 **MRSSing:** Eighteen professional singers participate in the recordings. Each singer is fluent in at
467 least one of the following languages: Chinese, English, German, or French, and collectively they
468 cover all vocal ranges including soprano, alto, tenor, and bass. Performers are paid \$50 per hour,
469 contributing a total of 75 hours of audio. Singing is recorded at a fixed position to ensure spatial
470 consistency. Binaural recordings are captured with head-mounted microphones, and monaural audio
471 is recorded with a studio-grade microphone. Third-person panoramic video is also captured.

472 **MRSMusic:** We engage 31 instrumentalists performing 23 Traditional Chinese, Western and Elec-
473 tronic instruments such as violin, erhu, pipa, electric guitar, and keyboard. Each performer receives
474 \$60 per recorded hour. A total of 75 hours of solo music performances are recorded. Audio is captured
475 using both head-mounted binaural microphones and reference monaural microphones. Third-person
476 video provides a full-scene view, while first-person cameras are used to capture detailed playing
477 gestures.

478 A.3 Details of Annotation

479 We employ a team of domain experts in singing, music performance, and linguistics to carry out and
480 review all annotations, compensating each annotator at a rate of \$15 per hour. Prior to beginning their
481 work, every expert receives a clear explanation of how the annotations will be used and agrees to
482 release their annotation results under an open-source license for academic research. For all modules,
483 we first synchronize each audio with its corresponding video, mono-channel reference track, and
484 3D positional metadata. Annotators then verify and correct this synchronization to ensure perfect
485 alignment across modalities.

486 **MRSLife.** In MRSDialogue scenes, we automatically generate initial transcripts and speaker clusters
487 with WhisperX, extracting word-level timestamps and speaker IDs. Experts then load these results in
488 Praat and assign each cluster to the correct speaker, correcting transcription errors as needed. Next,
489 we run Montreal Forced Aligner (MFA) (McAuliffe et al., 2017) to produce coarse phoneme-to-audio
490 alignments (exported in TextGrid format), using pypinyin to convert Chinese text into phoneme
491 sequences.⁵ Finally, annotators refine word and phoneme boundaries in Praat to achieve millisecond-
492 level precision. In MRSSound segments, annotators additionally label each time interval with the
493 corresponding event category (e.g., “clattering,” “typing,” “pages turning”).

494 **MRSSpeech.** Given full dialogue scripts, we perform automatic alignment to 30-minute recordings
495 using a chunk-based extension of WhisperX. This method divides each utterance-level audio segment
496 into fixed-length chunks (e.g., 30 seconds), applies phoneme-level models (e.g., wav2vec 2.0 (Baevski
497 et al., 2020)) for emission prediction on each chunk, and then concatenates emissions to form a
498 complete alignment matrix. We compute alignments via a trellis-based dynamic programming
499 algorithm with backtracking, followed by scaling to restore absolute timestamps. Word- and sentence-
500 level timings are derived by grouping aligned characters using word indices and sentence tokenization.
501 Missing or partial timings are interpolated using a specified method (e.g., nearest). This pipeline
502 enables efficient and accurate alignment of long-form recordings on GPU. We then refine sentence
503 boundaries in Praat and apply the same MFA-plus-Praat workflow used in MRSLife for fine-grained
504 phoneme and word-level alignment.

505 **MRSSing.** Each segment contains a solo vocal performance with known lyrics. We first apply
506 voice activity detection (VAD) to segment the recordings. Lyrics are aligned to each segment using
507 LyricFN’s ASR-based dynamic programming method. English phonemes follow the ARPABET
508 standard,⁶ while German and French use MFA’s phoneme sets.⁷ We then apply MFA for initial
509 alignment and refine it manually in Praat to obtain precise word and phoneme boundaries. Annotators
510 assign fine-grained style tags, including emotion, genre, and tempo. To generate score annotations,
511 we extract F0 contours using RMVPE (Wei et al., 2023) and convert them into MIDI format via
512 ROSVOT (Li et al., 2024), followed by expert correction.

⁵<https://github.com/mozillazg/python-pinyin>

⁶<https://en.wikipedia.org/wiki/ARPABET>

⁷<https://mfa-models.readthedocs.io/en/latest/dictionary/>

513 **MRSMusic.** We segment the recordings using AutoSlicer (Liu et al., 2022) and generate preliminary
514 symbolic annotations using basic-pitch (Bittner et al., 2022). Professional musicians then verify and
515 refine the annotations, adjusting note onsets, offsets, dynamics, and articulations to ensure consistency
516 between the score and performance.

517 **A.4 Details of Post-Processing**

518 **Segmentation:** After annotation, we segment the raw recordings into shorter clips to support spatial
519 audio generation and analysis tasks. For speech and singing, we use alignment timestamps to extract
520 utterances. For MRSSound, each recording is divided into fixed 10-second segments.

521 **Multimodal Synchronization:** In addition to binaural audio, each clip is synchronized with the
522 following modalities: audio, text, video, and position data. Using the segment timestamps, we align
523 all streams temporally. For static sources, we attach manually recorded 3D coordinates; for dynamic
524 sources, we interpolate Ultra-Wideband (UWB) tracking data over the segment duration. This process
525 yields fully synchronized multimodal clips ready for downstream spatial audio tasks. Furthermore,
526 for segments where participants’ faces are visible, we apply anonymization by using a face detection
527 model to overlay digital masks during post-processing.

528 **Checking:** To ensure the reliability of annotations, domain experts perform a random audit on
529 15% of the segmented clips across all four modules. The audit process involves the following
530 steps: (1) Verifying the temporal synchronization across different modal; (2) Confirming that the
531 assigned event labels accurately reflect the audiovisual content present in each clip; (3) Reviewing
532 the speech segments to check the correctness of word- and phoneme-level alignments; (4) Evaluating
533 singing clips for accurate alignment between lyrics and audio, consistency with musical scores,
534 and correctness of assigned style labels; (5) Assessing MRSMusic excerpts to verify that musical
535 properties, such as key, pitch, and note duration.

536 **A.5 Statistics of MRSAudio**

537 **A.5.1 Statistics of MRSLife**

538 **MRSDialogue.** Figure 4(a) presents the 3D spatial distribution of sound sources in MRSDialogue.
539 In this subset, which features frequent human conversations, most sources are located around the
540 ear-level height of the listener. The azimuthal distribution covers nearly all directions surrounding
541 the listener, offering diverse angular data for training spatial localization models with strong
542 generalization.

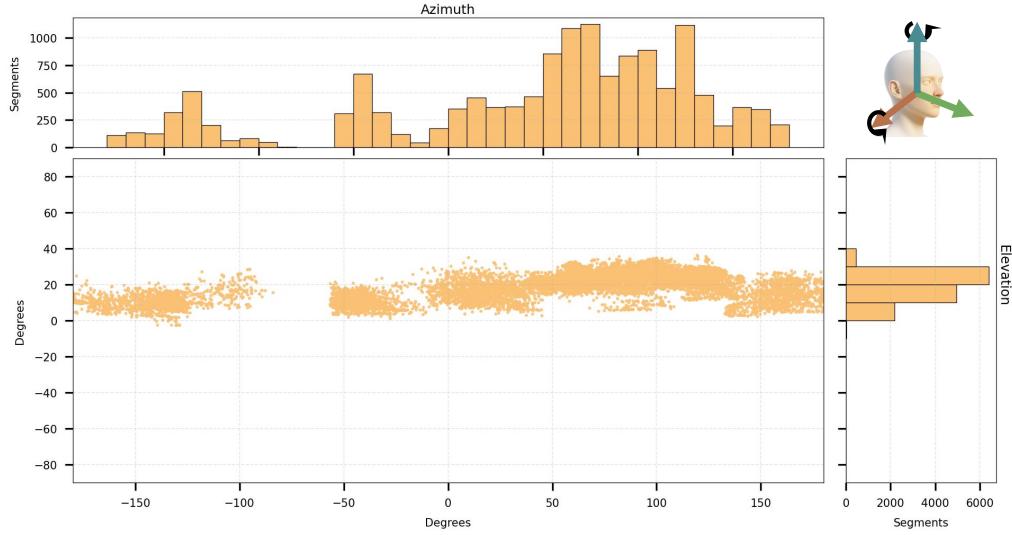
543 Figure 4(b) shows the phoneme distribution across all speech segments. The most common phoneme
544 is ‘i’, while the least frequent is ‘ueng’. This distribution aligns with real-world phonetic patterns
545 and highlights the dataset’s linguistic richness. Such broad phoneme coverage is beneficial for
546 downstream tasks in speech synthesis and recognition under spatial settings.

547 **MRSSound.** Figure 5(a) illustrates the spatial distribution of sound sources in MRSSound relative to
548 the listener’s head-centered coordinate system. The listener’s facing direction is at 90 degrees, and
549 most sound events are concentrated in the frontal hemisphere (azimuth from 0° to 180°), which is
550 consistent with the first-person video capture setup. Some events also appear in the rear field (-180°
551 to 0°). In elevation, the majority of sound sources are distributed between -90° and 40°, realistically
552 reflecting everyday human perception in standing scenarios, where sound events typically originate
553 below ear level. This broad spatial coverage supports the training of spatial audio models with strong
554 generalization ability.

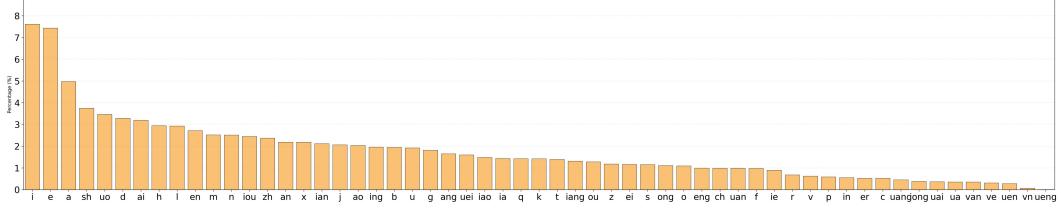
555 Figure 5(b) presents the duration distribution of recorded sound events. MRSSound covers a wide
556 variety of everyday scenarios, including cooking in kitchens, working in office environments, and
557 sports-related activities. The diversity of event types and durations makes the subset suitable for
558 training and evaluating models in real-world spatial sound event understanding.

559 **A.5.2 Statistics of MRSSpeech**

560 Figure 6(a) illustrates the spatial distribution of speech sources with respect to the listener. The
561 azimuth angles span the full 360° around the listener, providing comprehensive coverage of spatial
562 directions. Elevation angles are mostly concentrated between 0° and 60°. Smaller elevations reflect



(a) Spatial Distribution of Sound Sources Relative to the Listener



(b) Distribution of Phones in MRSDialogue

Figure 4: Statistical overview of MRSDialogue. (a) Spatial distribution of sound sources relative to the listener. Red, green, and blue arrows denote the positive x, y, and z axes; azimuth is measured around the z-axis from the x-axis, and elevation is relative to the xy-plane. (b) Distribution of phones in MRSDialogue.

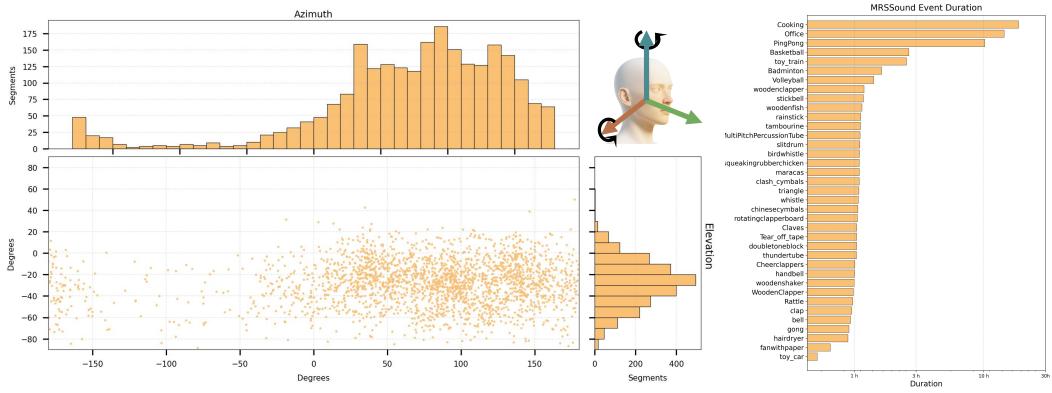


Figure 5: Statistical overview of MRSSound. (a) Spatial distribution of sound sources relative to the listener. Red, green, and blue arrows denote the positive x, y, and z axes; azimuth is measured around the z-axis from the x-axis, and elevation is relative to the xy-plane. (b) Duration distribution of sound event duration in MRSSound.

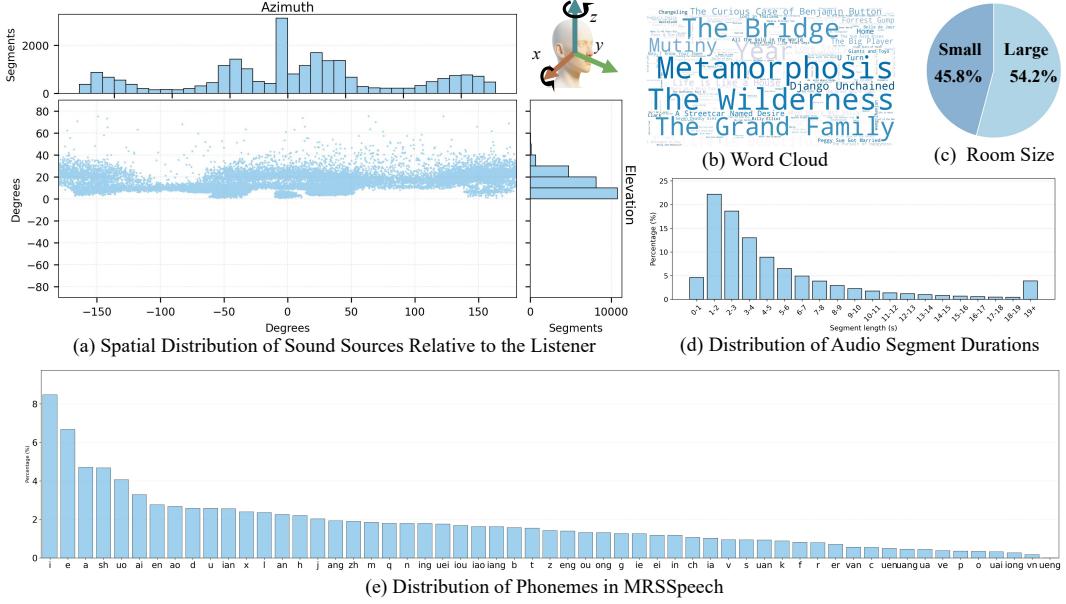


Figure 6: Statistical overview of MRSSpeech. (a) Spatial distribution of sound sources relative to the listener. Red, green, and blue arrows denote the positive x, y, and z axes; azimuth is measured around the z-axis from the x-axis, and elevation is relative to the xy-plane. (b) Word cloud. (c) Proportions of recording spaces by room size. (d) Distribution of audio segment durations. (e) Distribution of phonemes in MRSSpeech.

563 scenarios where both the speaker and listener are either standing or seated, while larger elevations
 564 (above 30°) simulate common real-world speech situations such as meetings, where a standing
 565 speaker addresses a seated listener. This diverse spatial coverage supports generalization in spatial
 566 speech modeling.

567 Figure 6(b) shows a word cloud representing the diversity of dialogue content. The transcripts are
 568 sourced from theatrical scripts, films, and other spoken-only scenarios, capturing a wide range of
 569 expressive and stylistic variation. Figure 6(c) presents the distribution of room sizes used for speech
 570 recordings. Most multi-speaker interactions take place in medium to large rooms, such as meeting
 571 or lecture spaces. We include three distinct environments with varying absorption properties and
 572 dimensions to simulate different acoustic conditions.

573 Figure 6(d) displays the distribution of audio segment durations. Most conversational turns are short,
 574 but the dataset also includes extended monologues exceeding 20 seconds, allowing models to capture
 575 both brief interactions and long-range motion or speaker dynamics.

576 Finally, Figure 6(e) presents the phoneme distribution, which covers all common phonetic units
 577 in the dataset’s target language. This phonetic diversity ensures that MRSSpeech provides strong
 578 generalizability for phoneme-aware models in speech synthesis and recognition.

579 A.5.3 Statistics of MRSSing

580 Figure 7(a) shows the 3D spatial distribution of sound sources with respect to the listener. The
 581 majority of sources are located in front of the listener, consistent with the setup of solo vocal
 582 recordings. However, the coverage also spans surrounding directions in both azimuth and elevation,
 583 ensuring spatial variability for training robust spatial audio models.

584 Figures 7 (b) and (c) highlight the diversity in style and content. Emotion annotations in Figure 7(b)
 585 include expressive labels such as happy and sad. Figure 7(c) displays the coverage across vocal
 586 ranges, including soprano, alto, tenor, and bass, ensuring a wide span of pitch and timbral variation.

587 Figure 7(d) presents the duration distribution of singing segments, which ranges from approximately
 588 4 to 10 seconds. This aligns with typical input lengths used in training singing voice synthesis

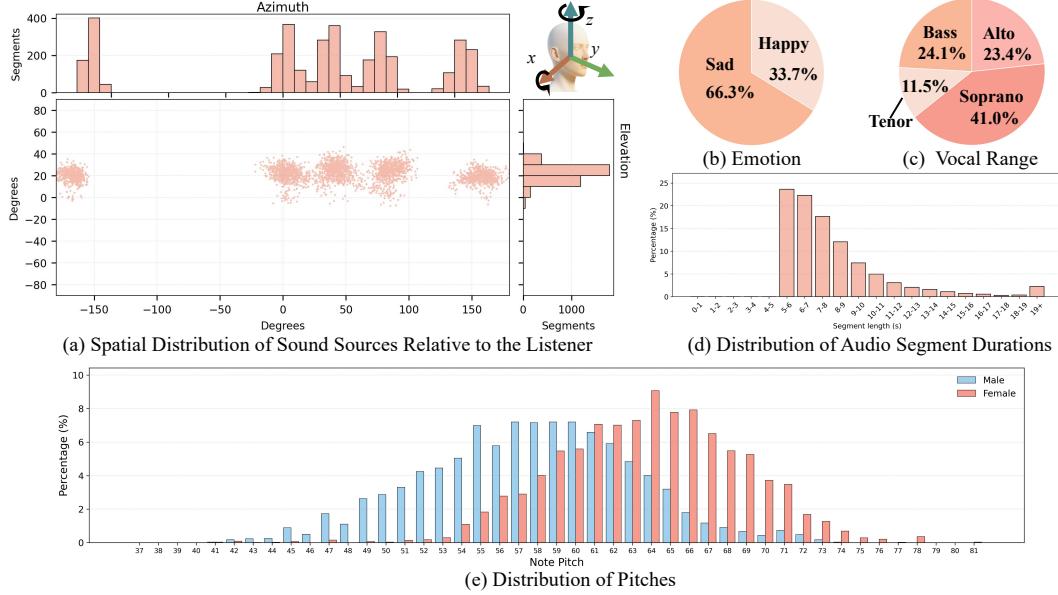


Figure 7: Statistical overview of MRSSing. (a) Spatial distribution of sound sources relative to the listener. Red, green, and blue arrows indicate the positive directions of the x, y, and z axes, respectively. Azimuth is measured around the z-axis from the x-axis, and elevation is relative to the xy-plane. (b) Emotion distribution. (c) Vocal range distribution. (d) Distribution of segment lengths. (e) Distribution of note pitches across segments.

models. Figure 7(e) illustrates the distribution of note pitches. The full range of musical notes is well represented, and a clear difference in pitch range is observed between male and female singers, with females generally singing at higher pitches. This confirms that MRSSing captures realistic vocal and musical variability.

In summary, MRSSing provides extensive diversity in spatial positioning, language, emotion, vocal range, segment duration, and pitch. This makes it a strong foundation for research on spatial singing voice synthesis and expressive vocal modeling.

596 A.5.4 Statistics of MRSMusic

597 Figure 8(a) illustrates the spatial positions of musical instruments relative to the listener. Most
598 sources are located in front of the listener, consistent with typical music listening scenarios. However,
599 the coverage also includes surrounding directions, contributing to spatial diversity in training and
600 evaluation.

601 Figures 8(b) and (c) highlight the diversity of instrument types and musical genres. The instrument
602 set spans Western, Traditional Chinese, and Electronic categories, while the genre annotations
603 include folk, pop, and classical music. Figure 8(f) further details the recording durations across
604 23 instruments, showing a relatively balanced distribution that facilitates downstream learning for
605 different instrument types.

606 Regarding generalization capacity, Figure 8(d) shows that audio segment durations range from
607 approximately 4 to 11 seconds, matching typical training input lengths. Figure 8(e) demonstrates
608 that the dataset covers a full range of musical pitch values, supporting tasks that require robust pitch
609 generalization.

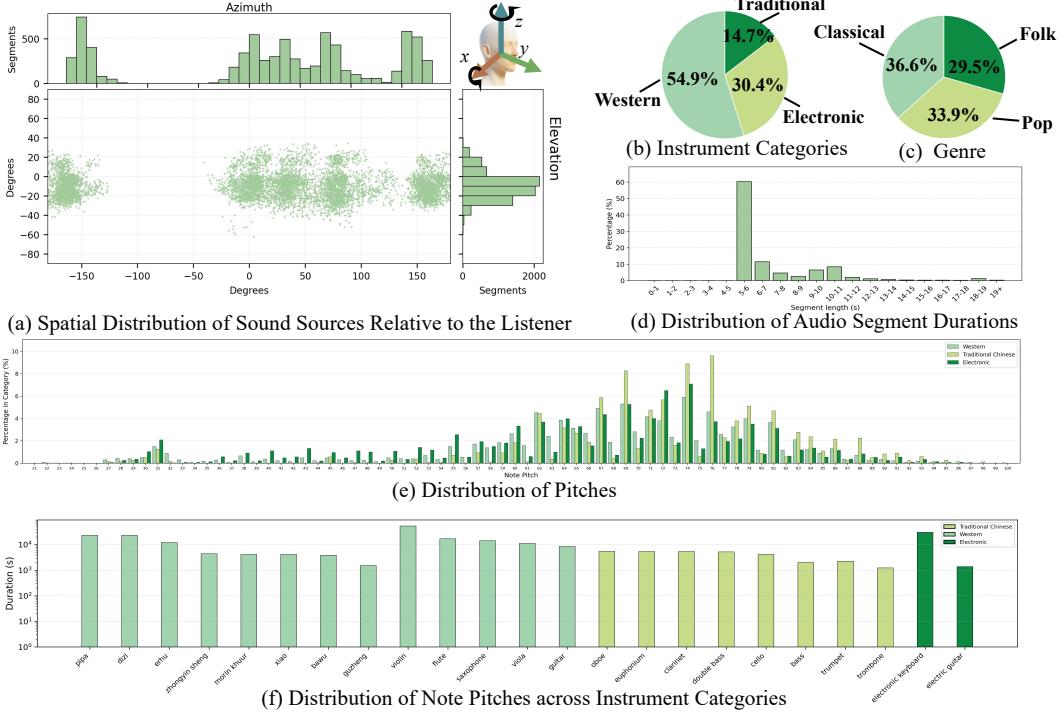


Figure 8: Statistical overview of MRSMusic. (a) Spatial distribution of sound sources relative to the listener. Red, green, and blue arrows denote the positive x, y, and z axes; azimuth is measured around the z-axis from the x-axis, and elevation is relative to the xy-plane. (b) Distribution of instrument categories duration. (c) Distribution of instrument genre duration. (d) Distribution of audio segment durations. (e) Distribution of pitches. (f) Distribution of instrument duration.

610 B Details of Experiments

611 B.1 Subjective Evaluation Metrics

612 We conduct subjective evaluation of the generation tasks using Mean Opinion Score (MOS). For
 613 each task, we randomly sample 40 utterances from the test set. Each utterance is paired with a
 614 corresponding source-position prompt and is evaluated by at least five expert listeners. MOS-Q,
 615 Listeners rate each sample on a five-point Likert scale ranging from 1 (bad) to 5 (excellent). For
 616 audio quality, we use M^OS-Q, where listeners wear headphones and assess the clarity and naturalness
 617 of the generated audio. For spatial perception, we use MOS-P, where listeners evaluate the realism of
 618 spatial cues and whether the perceived direction and distance of the sound source match the prompt
 619 description. All participants are fairly compensated for their time at a rate of \$15 per hour, resulting
 620 in a total cost of approximately \$2000. Participants are informed that their evaluations will be used
 621 exclusively for academic research purposes. Instructions for audio evaluations are shown in Figure 9.

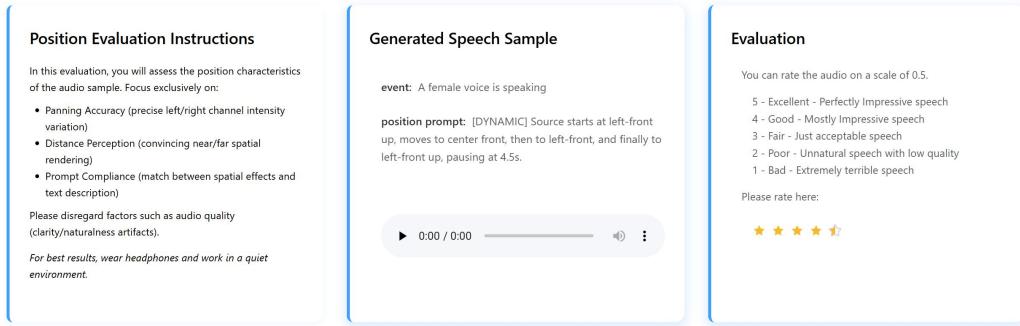
622 B.2 Objective Evaluation Metrics

623 To comprehensively assess spatial audio generation and understanding across multiple tasks, we
 624 adopt a set of objective metrics that evaluate signal fidelity, spatial consistency, intelligibility, and
 625 speaker or pitch accuracy. We randomly sample 400 data points as the test set.

626 **Audio Spatialization.** We measure waveform similarity using Wave L2, the mean squared error
 627 (MSE) between the generated and reference binaural waveforms. Amplitude L2 and Phase L2 are
 628 computed after applying Short-Time Fourier Transform (STFT), reflecting errors in the magnitude
 629 and phase components respectively. MRSTFT loss⁸ is also used, combining spectral convergence,

⁸<https://github.com/csteinmetz1/auraloss>

MOS-P Testing



(a) Screenshot of MOS-P Testing

MOS-Q Testing



(b) Screenshot of MOS-Q Testing

Figure 9: Instructions for audio evaluations. (a) Screenshot of MOS-P Testing. (b) Screenshot of MOS-Q Testing

630 log-magnitude, and linear-magnitude terms for better spectral alignment. In addition, we use the
 631 perceptual evaluation metric PESQ⁹ to assess audio quality for speech-related tasks. Since PESQ is
 632 designed for speech, it is omitted for MRSLife and MRSMusic. For all metrics except PESQ, lower
 633 values indicate better performance.

634 To evaluate spatial consistency, we use Spatial-AST(Zheng et al., 2024) to extract angular and distance
 635 embeddings from the binaural audio. Since Spatial-AST predicts positions only for static sources, we
 636 compute the cosine similarity between the predicted and ground truth embeddings within 1-second
 637 segments and average the results to assess overall spatial fidelity.

638 **Spatial Text to Speech.** We evaluate speech intelligibility using Character Error Rate (CER), which
 639 measures the proportion of character-level differences between ASR transcriptions and reference
 640 texts. Transcriptions are generated using the Paraformer-zh model (Gao et al., 2023). To assess
 641 speaker consistency, we compute Speaker Identity Matching (SIM) as the cosine similarity between
 642 speaker embeddings extracted using a WavLM-based speaker verification model.¹⁰

643 **Spatial Singing Voice Synthesis.** We use Mel Cepstral Distortion (MCD) to assess the spectral
 644 similarity between the generated and reference vocals. It is defined as:

⁹<https://github.com/aliutkus/speechmetrics>

¹⁰<https://huggingface.co/pyannote/speaker-diarization>

$$\text{MCD} = \frac{10}{\ln 10} \sqrt{2 \sum_{d=1}^D (m_t(d) - \hat{m}_t(d))^2}, \quad (1)$$

645 where $m_t(d)$ and $\hat{m}_t(d)$ represent the d -th Mel-frequency cepstral coefficient (MFCC) at frame t for
 646 the ground truth and synthesized signals, respectively, and D is the number of MFCC dimensions.

647 Additionally, we adopt F0 Frame Error (FFE) to evaluate pitch accuracy by comparing extracted F0
 648 contours between the synthesized and ground truth audio.

649 **Spatial Music Generation.** We use Fréchet Audio Distance (FAD) (Kilgour et al., 2018) to assess
 650 perceptual similarity between the feature distributions of generated and reference audio. In addition,
 651 F0 Frame Error (FFE) is used to evaluate pitch accuracy by comparing the extracted F0 contours
 652 against the reference musical scores.

653 **Sound Event Localization and Detection (SELD).** Following STARSS23 (Shimada et al., 2023),
 654 we adopt four joint detection and localization metrics: F_{20° : location-aware F-score; a prediction is
 655 correct if the event class matches and angular error is below 20° . ER_{20° : error rate computed as the
 656 sum of insertions, deletions, and substitutions over reference events. LE_{CD} : class-aware localization
 657 error, the mean angular difference between predicted and reference directions. LR_{CD} : class-aware
 658 localization recall, the percentage of correctly localized events among all instances of each class. We
 659 compute all metrics in 1-second non-overlapping segments using macro-averaging across all event
 660 classes. Higher values of F_{20° and LR_{CD} , and lower values of ER_{20° and LE_{CD} indicate better
 661 performance.

662 B.3 Audio Spatialization

663 We build upon BinauralGrad’s two-stage diffusion-based framework to convert monaural audio into
 664 spatial audio. However, in our subsequent generation experiments, the synthesized monaural audio
 665 typically corresponds to a centrally positioned source between the ears, effectively serving as the
 666 first stage of BinauralGrad. Therefore, for the spatialization experiments presented here, we use only
 667 the second stage of BinauralGrad to validate spatial audio generation. Specifically, we first convert
 668 binaural recordings to monaural input by averaging the two channels. Next, we apply a DSP-based
 669 method to produce a coarse spatial approximation of the binaural signal. This monaural input and its
 670 simulated binaural counterpart are then used as input to the BinauralGrad model, which is conditioned
 671 on the object’s motion trajectory to generate spatialized binaural audio. We list the architecture and
 672 hyperparameters of BinauralGrad in Table 8.

Table 8: Hyper-parameters of BinauralGrad modules.

Hyperparameter		BinauralGrad
Binaural Encoder	Wave Encoder Layers	2
	Position Encoder Layers	2
	Encoder Conv1D Kernel	3
	Encoder Dropout	0.4
Mel Predictor	Residual Blocks	3
	Bidirectional Layers	3
	Hidden_size	128
	Training Steps	200
	Sampling Steps	6

673 B.4 Spatial Text to Speech

674 We fine-tune two pre-trained models, CosyVoice and F5-TTS, using monaural audio obtained by
 675 averaging the binaural recordings from the MRSSpeech subset. This allows the models to learn the
 676 generation characteristics specific to our dataset. After monaural generation, we apply a pre-trained
 677 spatialization model trained on MRSSpeech to convert the output into spatialized binaural audio.

678 For the ISDrama baseline, we adopt the model variant proposed in the original paper that conditions
 679 generation on predefined spatial paths to produce binaural spatial speech directly.

Table 9: Hyper-parameters of Rmssinger modules.

Hyperparameter		Rmssinger
Phoneme Encoder	Phoneme Embedding	256
	Encoder Layers	4
	Encoder Hidden	256
	Encoder Conv1D Kernel	9
	Encoder Conv1D Filter Size	1024
	Encoder Attention Heads	2
Note Encoder	Encoder Dropout	0.1
	Pitches Embedding	256
	Type Embedding	256
Pitch Predictor	Duration Hidden	256
	Conv Layers	12
	Kernel Size	3
	Residual Channel	192
	Hidden Channel	256
Mel Predictor	Training Steps	100
	Conv Layers	20
	Kernel Size	3
	Residual Channel	256
	Hidden Channel	256
	Training Steps	100

680 B.5 Spatial Singing Voice Synthesis

681 For the ISDrama baseline, we adopt the spatial-path-conditioned generation framework and extend it
 682 by incorporating Rmssinger’s note encoder, allowing explicit pitch control through musical score
 683 input. We use the single stage variation of Rmssinger, a standard singing voice synthesis model, and
 684 train it with monaural targets derived from the averaged binaural recordings in the MRSSing subset.
 685 This enables the model to generate pitch-accurate audio that reflects the acoustic characteristics. The
 686 synthesized monaural audio is then spatialized using the spatialization model trained on MRSSing.
 687 We list the architecture and hyperparameters of Rmssinger in Table 9.

688 B.6 Spatial Music Generation

689 For spatial music generation, in the ISDrama setting, we follow the path-conditioned generation
 690 pipeline but remove the phoneme encoder, using only the note encoder to process symbolic scores
 691 as the sole control condition for generating spatialized music. We also adopt Make-An-Audio 2
 692 as our base model and train it using monaural audio derived from averaged MRSMusic binaural
 693 recordings. We use the pretrained spectrogram autoencoder. Instrument category and symbolic
 694 score information are concatenated into the text prompt to guide the music generation process. The
 695 generated audio is subsequently spatialized using the MRSMusic-trained spatialization model. We
 696 list the hyper-parameters of Make-An-Audio 2 in Table 10.

697 B.7 Sound Event Localization and Detection

698 We follow the STARSS23 framework for FOA-based sound event localization and detection, training
 699 on the event segments from the MRSSound subset under both audio-only and audio-visual conditions.
 700 To enable binaural input, we modify the model architecture to use three input channels and extract
 701 interaural phase difference features. This allows us to adapt the SELD model to binaural audio.
 702 Additionally, we experiment with replacing convolutional layers with Transformer encoders to explore
 703 the effect of different architectures on sound event localization and detection performance. We list
 704 the hyper-parameters of SELD in Table 11.

Table 10: Hyperparameters of Make-An-Audio 2.

Hyperparameter		Make-An-Audio 2
Autoencoders	Input/Output Channels	80
	Hidden Channels	20
	Residual Blocks	2
	Spectrogram Shape	(80, 624)
	Channel Multipiler	[1, 2, 4]
Transformer Backbone	Input shape	(20, T)
	Condition_embedding Size	1024
	Feed-forward Hidden_size	576
	Num of Transformer Heads	8
	Transformer Blocks	8
	Training Steps	1000
	Sampling Steps	100
CLAP Text Encoder	Transformer Embed Channels	768
	Output Project Channels	1024
	Token Length	77

Table 11: Hyperparameters of SELD.

Hyperparameter		SELD
Input	Binaural Channels	3
	FOA Channels	7
	Frequency Bins	128
	Frames	120
Audio Encoder(CNN)	Hidden_size	64
	Conv Blocks	3
Audio Encoder(Vit)	Hidden_size	128
	Conv Blocks	1
	Transformer Blocks	2
	Num of Transformer Heads	4

705 **C Authorship Statement**

706 **Zhou Zhao** and **Fei Wu** served as the overall project leaders and supervisors, offering full-scale
707 resource support throughout the development of the MRSAudio project.

708 **Data Planning** **Wenxiang Guo** is responsible for all aspects of the MRSAudio project, from its
709 initial conception to its final execution. The planning and design of the individual sub-datasets are
710 undertaken by:

- 711 • **MRSSpeech**: Yu Zhang, Zhiyuan Zhu, Wenxiang Guo
- 712 • **MRSSing**: Changhao Pan, Yu Zhang
- 713 • **MRSMusic**: Changhao Pan, Yu Zhang
- 714 • **MRSAudio**: Wenxiang Guo, Zhiyuan Zhu, Li Tang

715 **Data Recording**

- 716 • **MRSSpeech**: **Xintong Hu** serves as the principal coordinator of the dataset, taking charge of
717 venue management, volunteer recruitment, audio recording, and data organization. **Yixuan**
718 **Chen** and **Pengfei Fan** assist in the recording and organization of audio data, as well as the
719 setup of recording environments. **Zhetao Chen**, **Yanhao Yu**, **Ke Xu**, and **Qiang Huang**
720 participate in the audio recording and data organization. **Wenxiang Guo** provides technical
721 support during the recording process.
- 722 • **MRSSing**: **Changhao Pan** serves as the recording lead for this sub-dataset, taking charge of
723 venue management, volunteer recruitment, audio recording, and data organization. **Yuhan**
724 **Wang** assists with the recruitment and coordination of recording personnel and primarily
725 handles data organization. **Ke Xu** and **Li Tang** participate in the recording and organization
726 of audio data.
- 727 • **MRSMusic**: **Changhao Pan** serves as the recording lead for this sub-dataset. **Changhao**
728 **Pan** and **Hankun Xu** are responsible for the recruitment and coordination of performers.
729 **Han Wang**, **Hankun Xu**, and **Changhao Pan** jointly conduct the audio recording, with
730 **Han Wang** primarily responsible for data organization.
- 731 • **MRSLife**: **Wenxiang Guo** serves as the lead coordinator for the MRSLife sub-dataset,
732 which consists of two parts: *MRSSound* and *MRSDialogue*. **Wenxiang Guo**, **Rui Yang**,
733 **Zhiyuan Zhu**, and **Xintong Hu** jointly carry out the recording of MRSSound. **Zhiyuan**
734 **Zhu** is responsible for the scene design of MRSDialogue, while **Yuhan Wang**, **Rui Yang**,
735 and **Zhiyuan Zhu** conduct the recording of MRSDialogue.

736 **Data Processing and Annotation**

- 737 • **MRSSpeech**: **Zhiyuan Zhu** is responsible for the complete data processing pipeline, includ-
738 ing segmentation, ASR, alignment, and the organization of spatial annotation. **Wenxiang**
739 **Guo** provides technical support and guidance throughout the process.
- 740 • **MRSSing**: **Changhao Pan** and **Zongbao Zhang** process the dataset, covering segmentation,
741 text alignment, note transcription, and the alignment of spatial information. **Yuhan Wang** is
742 responsible for global style annotation, while **Han Wang** organizes the music score data.
- 743 • **MRSMusic**: **Changhao Pan** and **Zongbao Zhang** handle data processing, including seg-
744 mentation, note transcription, and spatial information alignment. **Hankun Xu** is responsible
745 for both global style annotation and music score organization.
- 746 • **MRSLife**: **Zhiyuan Zhu** and **Wenxiang Guo** jointly perform the data processing. **Rui**
747 **Yang** annotates the audio events in the recordings.

748 **Experiments and Benchmarks**

- 749 • **Sound Event Localization and Detection (SELD)**: **Wenxiang Guo** updates and adapts the
750 baseline models for this task. **Zhiyuan Zhu** is responsible for training and evaluating the
751 baseline models. **Changhao Pan** provides technical guidance for the reproduction process.

- 752 • **Audio Spatialization:** **Wenxiang Guo** is solely responsible for the design, implementation,
753 and evaluation of this benchmark.
- 754 • **Text-to-Speech:** **Yu Zhang** conducts the experimental work for this benchmark, while
755 **Wenxiang Guo** performs model evaluation.
- 756 • **Singing Voice Synthesis:** **Changhao Pan** and **Yu Zhang** carry out the experiments for this
757 benchmark. **Xintong Hu** is responsible for model evaluation.
- 758 • **Spatial Music Generation:** **Changhao Pan**, **Wenxiang Guo**, and **Yu Zhang** jointly conduct
759 the experiments and evaluations for this benchmark.

760 **Paper Writing** **Wenxiang Guo** is responsible for the overall writing of the paper. **Changhao Pan**
761 contributes to the writing of the experimental sections. **Zhiyuan Zhu** is in charge of the writing of
762 the statistical analysis section. **Xintong Hu**, **Yu Zhang**, and **Zhou Zhao** participate in revising and
763 polishing the manuscript.

764 **NeurIPS Paper Checklist**

765 **1. Claims**

766 Question: Do the main claims made in the abstract and introduction accurately reflect the
767 paper's contributions and scope?

768 Answer: [Yes]

769 Justification: As shown in Section 1, this paper proposes a large-scale multimodal recorded
770 spatial audio dataset along with refined annotation for spatial audio-related tasks.

771 Guidelines:

- 772 • The answer NA means that the abstract and introduction do not include the claims
773 made in the paper.
- 774 • The abstract and/or introduction should clearly state the claims made, including the
775 contributions made in the paper and important assumptions and limitations. A No or
776 NA answer to this question will not be perceived well by the reviewers.
- 777 • The claims made should match theoretical and experimental results, and reflect how
778 much the results can be expected to generalize to other settings.
- 779 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
780 are not attained by the paper.

781 **2. Limitations**

782 Question: Does the paper discuss the limitations of the work performed by the authors?

783 Answer: [Yes]

784 Justification: We have discussed the limitations and the negative societal impact in Section 5.

785 Guidelines:

- 786 • The answer NA means that the paper has no limitation while the answer No means that
787 the paper has limitations, but those are not discussed in the paper.
- 788 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 789 • The paper should point out any strong assumptions and how robust the results are to
790 violations of these assumptions (e.g., independence assumptions, noiseless settings,
791 model well-specification, asymptotic approximations only holding locally). The authors
792 should reflect on how these assumptions might be violated in practice and what the
793 implications would be.
- 794 • The authors should reflect on the scope of the claims made, e.g., if the approach was
795 only tested on a few datasets or with a few runs. In general, empirical results often
796 depend on implicit assumptions, which should be articulated.
- 797 • The authors should reflect on the factors that influence the performance of the approach.
798 For example, a facial recognition algorithm may perform poorly when image resolution
799 is low or images are taken in low lighting. Or a speech-to-text system might not be
800 used reliably to provide closed captions for online lectures because it fails to handle
801 technical jargon.
- 802 • The authors should discuss the computational efficiency of the proposed algorithms
803 and how they scale with dataset size.
- 804 • If applicable, the authors should discuss possible limitations of their approach to
805 address problems of privacy and fairness.
- 806 • While the authors might fear that complete honesty about limitations might be used by
807 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
808 limitations that aren't acknowledged in the paper. The authors should use their best
809 judgment and recognize that individual actions in favor of transparency play an impor-
810 tant role in developing norms that preserve the integrity of the community. Reviewers
811 will be specifically instructed to not penalize honesty concerning limitations.

812 **3. Theory assumptions and proofs**

813 Question: For each theoretical result, does the paper provide the full set of assumptions and
814 a complete (and correct) proof?

815 Answer: [NA]

816 Justification: This paper does not include theoretical results.
817 Guidelines:

- 818 • The answer NA means that the paper does not include theoretical results.
- 819 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 820 • referenced.
- 821 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 822 • The proofs can either appear in the main paper or the supplemental material, but if
- 823 • they appear in the supplemental material, the authors are encouraged to provide a short
- 824 • proof sketch to provide intuition.
- 825 • Inversely, any informal proof provided in the core of the paper should be complemented
- 826 • by formal proofs provided in appendix or supplemental material.
- 827 • Theorems and Lemmas that the proof relies upon should be properly referenced.

828 **4. Experimental result reproducibility**

829 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
830 perimental results of the paper to the extent that it affects the main claims and/or conclusions
831 of the paper (regardless of whether the code and data are provided or not)?

832 Answer: [Yes]

833 Justification: Our dataset is available at the <https://huggingface.co/datasets/verstar/MRSAudio>. The source code associated with our work can be accessed in the GitHub repository: https://github.com/MRSAudio/MRSAudio_Main and in the supplementary materials. Additionally, training details is shown in Appendix B.

837 Guidelines:

- 838 • The answer NA means that the paper does not include experiments.
- 839 • If the paper includes experiments, a No answer to this question will not be perceived
- 840 • well by the reviewers: Making the paper reproducible is important, regardless of
- 841 • whether the code and data are provided or not.
- 842 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 843 • to make their results reproducible or verifiable.
- 844 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 845 • For example, if the contribution is a novel architecture, describing the architecture fully
- 846 • might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 847 • be necessary to either make it possible for others to replicate the model with the same
- 848 • dataset, or provide access to the model. In general, releasing code and data is often
- 849 • one good way to accomplish this, but reproducibility can also be provided via detailed
- 850 • instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 851 • of a large language model), releasing of a model checkpoint, or other means that are
- 852 • appropriate to the research performed.
- 853 • While NeurIPS does not require releasing code, the conference does require all submissions
- 854 • to provide some reasonable avenue for reproducibility, which may depend on the
- 855 • nature of the contribution. For example
 - 856 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
 - 857 • to reproduce that algorithm.
 - 858 (b) If the contribution is primarily a new model architecture, the paper should describe
 - 859 • the architecture clearly and fully.
 - 860 (c) If the contribution is a new model (e.g., a large language model), then there should
 - 861 • either be a way to access this model for reproducing the results or a way to reproduce
 - 862 • the model (e.g., with an open-source dataset or instructions for how to construct
 - 863 • the dataset).
 - 864 (d) We recognize that reproducibility may be tricky in some cases, in which case
 - 865 • authors are welcome to describe the particular way they provide for reproducibility.
 - 866 • In the case of closed-source models, it may be that access to the model is limited in
 - 867 • some way (e.g., to registered users), but it should be possible for other researchers
 - 868 • to have some path to reproducing or verifying the results.

869 **5. Open access to data and code**

870 Question: Does the paper provide open access to the data and code, with sufficient instruc-
871 tions to faithfully reproduce the main experimental results, as described in supplemental
872 material?

873 Answer: [Yes]

874 Justification: Our dataset is available at the [https://huggingface.co/datasets/](https://huggingface.co/datasets/verstar/MRSAudio)
875 [verstar/MRSAudio](https://github.com/MRSAudio/MRSAudio_Main). The source code associated with our work can be accessed in the
876 GitHub repository: https://github.com/MRSAudio/MRSAudio_Main and in the sup-
877 plementary materials.

878 Guidelines:

- 879 • The answer NA means that paper does not include experiments requiring code.
- 880 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/](https://nips.cc/public/guides/CodeSubmissionPolicy)
881 [public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 882 • While we encourage the release of code and data, we understand that this might not be
883 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
884 including code, unless this is central to the contribution (e.g., for a new open-source
885 benchmark).
- 886 • The instructions should contain the exact command and environment needed to run to
887 reproduce the results. See the NeurIPS code and data submission guidelines ([https://nips.cc/](https://nips.cc/public/guides/CodeSubmissionPolicy)
888 [public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 889 • The authors should provide instructions on data access and preparation, including how
890 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 891 • The authors should provide scripts to reproduce all experimental results for the new
892 proposed method and baselines. If only a subset of experiments are reproducible, they
893 should state which ones are omitted from the script and why.
- 894 • At submission time, to preserve anonymity, the authors should release anonymized
895 versions (if applicable).
- 896 • Providing as much information as possible in supplemental material (appended to the
897 paper) is recommended, but including URLs to data and code is permitted.

898 6. Experimental setting/details

899 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
900 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
901 results?

902 Answer: [Yes]

903 Justification: The settings and details of experiments is shown in Appendix B.

904 Guidelines:

- 905 • The answer NA means that the paper does not include experiments.
- 906 • The experimental setting should be presented in the core of the paper to a level of detail
907 that is necessary to appreciate the results and make sense of them.
- 908 • The full details can be provided either with the code, in appendix, or as supplemental
909 material.

910 7. Experiment statistical significance

911 Question: Does the paper report error bars suitably and correctly defined or other appropriate
912 information about the statistical significance of the experiments?

913 Answer: [Yes]

914 Justification: We report confidence intervals of subjective metric results in Section 4.

915 Guidelines:

- 916 • The answer NA means that the paper does not include experiments.
- 917 • The authors should answer “Yes” if the results are accompanied by error bars, confi-
918 dence intervals, or statistical significance tests, at least for the experiments that support
919 the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The required experimental resources are detailed in Section 4 and Appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: 1) The wage details for participants are reported in Appendix A.2, A.3, and B.1. 2) All data were de-identified prior to release, as detailed in Section 3. 3) As described in Appendix A.2, all participants involved in data recording signed consent forms, agreeing to the publication of the dataset under the CC BY-NC-SA 4.0 license. 4) The dataset we have released adheres to the CC BY-NC-SA 4.0 license, and all baselines used in the paper comply with their respective licenses.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

971 Justification: We discuss the societal impacts from both positive and negative perspectives
972 in Section 5.

973 Guidelines:

- 974 • The answer NA means that there is no societal impact of the work performed.
- 975 • If the authors answer NA or No, they should explain why their work has no societal
976 impact or why the paper does not address societal impact.
- 977 • Examples of negative societal impacts include potential malicious or unintended uses
978 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
979 (e.g., deployment of technologies that could make decisions that unfairly impact specific
980 groups), privacy considerations, and security considerations.
- 981 • The conference expects that many papers will be foundational research and not tied
982 to particular applications, let alone deployments. However, if there is a direct path to
983 any negative applications, the authors should point it out. For example, it is legitimate
984 to point out that an improvement in the quality of generative models could be used to
985 generate deepfakes for disinformation. On the other hand, it is not needed to point out
986 that a generic algorithm for optimizing neural networks could enable people to train
987 models that generate Deepfakes faster.
- 988 • The authors should consider possible harms that could arise when the technology is
989 being used as intended and functioning correctly, harms that could arise when the
990 technology is being used as intended but gives incorrect results, and harms following
991 from (intentional or unintentional) misuse of the technology.
- 992 • If there are negative societal impacts, the authors could also discuss possible mitigation
993 strategies (e.g., gated release of models, providing defenses in addition to attacks,
994 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
995 feedback over time, improving the efficiency and accessibility of ML).

996 11. Safeguards

997 Question: Does the paper describe safeguards that have been put in place for responsible
998 release of data or models that have a high risk for misuse (e.g., pretrained language models,
999 image generators, or scraped datasets)?

1000 Answer: [Yes]

1001 Justification: All data were de-identified prior to release, as detailed in Section 3 and
1002 Appendix A.4. For participants whose faces were visible, we anonymized the recordings by
1003 adding masks.

1004 Guidelines:

- 1005 • The answer NA means that the paper poses no such risks.
- 1006 • Released models that have a high risk for misuse or dual-use should be released with
1007 necessary safeguards to allow for controlled use of the model, for example by requiring
1008 that users adhere to usage guidelines or restrictions to access the model or implementing
1009 safety filters.
- 1010 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1011 should describe how they avoided releasing unsafe images.
- 1012 • We recognize that providing effective safeguards is challenging, and many papers do
1013 not require this, but we encourage authors to take this into account and make a best
1014 faith effort.

1015 12. Licenses for existing assets

1016 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1017 the paper, properly credited and are the license and terms of use explicitly mentioned and
1018 properly respected?

1019 Answer: [Yes]

1020 Justification: We use all models and codes under their Licenses (e.g. CC 4.0, MIT, etc).

1021 Guidelines:

- 1022 • The answer NA means that the paper does not use existing assets.
- 1023 • The authors should cite the original paper that produced the code package or dataset.

- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We provide detailed description in Section 3 and Appendix A for the proposed dataset.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [Yes]

Justification: Appendix B.1 includes the instructions and screenshots of rating systems. The wage for participants is mentioned in Appendix B.1 as well.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

1076 Guidelines:

- 1077 • The answer NA means that the paper does not involve crowdsourcing nor research with
1078 human subjects.
- 1079 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1080 may be required for any human subjects research. If you obtained IRB approval, you
1081 should clearly state this in the paper.
- 1082 • We recognize that the procedures for this may vary significantly between institutions
1083 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1084 guidelines for their institution.
- 1085 • For initial submissions, do not include any information that would break anonymity (if
1086 applicable), such as the institution conducting the review.

1087 **16. Declaration of LLM usage**

1088 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1089 non-standard component of the core methods in this research? Note that if the LLM is used
1090 only for writing, editing, or formatting purposes and does not impact the core methodology,
1091 scientific rigorousness, or originality of the research, declaration is not required.

1092 Answer: [NA]

1093 Justification: The core method development in this research does not involve LLMs as
1094 any important, original, or non-standard components. 1) About LLM: this work only uses
1095 pretrained language models to generate natural prompts based on the given trajectory as
1096 shown in Section 3; 2) We only use image generators to assist the design of the paper's
1097 header figure.

1098 Guidelines:

- 1099 • The answer NA means that the core method development in this research does not
1100 involve LLMs as any important, original, or non-standard components.
- 1101 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
1102 for what should or should not be described.