

## 1 **A Supplementary materials**

### 2 **A.1 Documentation and intended uses**

3 We include a datasheet in Section B.

### 4 **A.2 Data access and code availability**

5 We have released the RealMAN dataset on the AISHELL website to ensure its long-term  
6 preservation. Researchers can access our dataset either directly at the dataset website <https://www.aishelltech.com/RealMAN> or find the dataset website at the GitHub repository <https://github.com/Audio-WestlakeU/RealMAN>.  
7 //www.aishelltech.com/RealMAN or find the dataset website at the GitHub repository <https://github.com/Audio-WestlakeU/RealMAN>.  
8 //github.com/Audio-WestlakeU/RealMAN.

9 The detailed information and code about the RealMAN dataset can be found at the GitHub repository  
10 <https://github.com/Audio-WestlakeU/RealMAN>. This link also gives a detailed description  
11 of how to use recordings from the RealMAN dataset for speech enhancement and localization.

### 12 **A.3 Data format**

13 RealMAN dataset provides audio files in the format of ‘flac’, position annotation files in the format  
14 of ‘csv’, and a transcription file in the plain text format of ‘trn’, all of which can be read by a variety  
15 of tools.

### 16 **A.4 Hosting and maintenance plan**

17 RealMAN will remain hosted on the AISHELL website and maintained by the AudioLab of Westlake  
18 University for the foreseeable future. Any changes will be updated on the Github repository.

### 19 **A.5 Statement of responsibility**

20 The authors declare that they bear all responsibility for violations of rights and that this dataset  
21 is released under CC-BY-4.0 license. During the recording process, we adhered to the Personal  
22 Information Protection Law of the People’s Republic of China (China’s PIPL) and explained all  
23 privacy-related details to the participants. We ensured the participants were fully informed and  
24 obtained the consent of all involved.

### 25 **A.6 License**

26 The RealMAN dataset is made available under the CC BY 4.0 license. The authors bear all responsi-  
27 bility in case of violation of rights.

### 28 **A.7 Potential negative societal impacts**

29 Our dataset has no potential negative societal impact.

### 30 **A.8 Personally identifiable information or offensive content**

31 Our dataset has no personally identifiable information or offensive content.

## 32 **B Datasheet**

### 33 **B.1 Motivation**

#### 34 **1. For what purpose was the dataset created?**

35 The training of deep learning-based multichannel speech enhancement and source localization systems  
36 relies heavily on the simulation of room impulse response and multichannel diffuse noise, due to the

37 lack of large-scale real-recorded datasets. However, the acoustic mismatch between simulated and  
38 real-world data could degrade the model performance when applying in real-world scenarios. To  
39 bridge this simulation-to-real gap, this paper presents a new relatively large-scale Real-recorded and  
40 annotated Microphone Array speech&Noise (RealMAN) dataset.

41 **2. Who created this dataset (e.g., which team, research group) and on behalf of which entity**  
42 **(e.g., company, institution, organization)**

43 The RealMAN dataset was created by AudioLab of Westlake University and Beijing AIShell Tech-  
44 nology Co. Ltd.

45 **3. Who funded the creation of the dataset?**

46 This work was supported by the Zhejiang Provincial Natural Science Foundation of China under  
47 Grant 2022XHSJJ008 and the Postdoctoral Science Foundation of China under Grant 2022M722848.

## 48 **B.2 Composition**

49 **1. What do the instances that comprise the dataset represent (e.g., documents, photos, people,**  
50 **countries)?**

51 The instances of RealMAN are a series of files, which include real-recorded 32-channel speech and  
52 noise recordings in the format of ‘flac’, direct-path clean signals in the format of ‘flac’, and position  
53 annotation files in the format of ‘csv’.

54 **2. How many instances are there in total (of each type, if appropriate)?**

55 The dataset consists of 83 hours of speech and 144 hours of noise, recorded in 32 and 31 different  
56 scenes, respectively. In addition, this dataset provides annotations of source azimuth angle, direct-path  
57 target clean speech, and speech transcription.

58 **3. Does the dataset contain all possible instances or is it a sample (not necessarily random) of**  
59 **instances from a larger set?**

60 The RealMAN dataset is recorded from scratch, and is not sampled from a larger set.

61 **4. What data does each instance consist of?**

62 Each instance of the RealMAN dataset consists of microphone signal recordings and the correspond-  
63 ing annotations.

64 **5. Is there a label or target associated with each instance?**

65 Yes. The RealMAN provides a series of labels/targets including source azimuth angle, direct-path  
66 target clean speech, and speech transcription, for speech enhancement and source localization.

67 **6. Is any information missing from individual instances?**

68 No.

69 **7. Are relationships between individual instances made explicit (e.g., users’ movie ratings, social**  
70 **network links)?**

71 N/A

72 **8. Are there recommended data splits (e.g., training, development/validation, testing)?**

73 Yes. We split them into training, validation and testing sets according to the acoustic characteristics  
74 of the recording scenes and speaker identities. The total 83.9 hours of the recorded speech are divided  
75 into 63.5, 7.8 and 11.6 hours for training, validation and test, respectively. And 144.5 hours of noise  
76 data are divided into 106.3, 16.0 and 22.2 hours for training, validation and test, respectively.

77 **9. Are there any errors, sources of noise, or redundancies in the dataset?**

78 No.

79 **10. Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,**  
80 **websites, tweets, other datasets)?**

81 The RealMAN dataset is mostly self-contained, but the speech source signals played by the loud-  
82 speaker are obtained from AISHELL.

83 **11. Does the dataset contain data that might be considered confidential (e.g., data that is**  
84 **protected by legal privilege or by doctor-patient confidentiality, data that includes the content**  
85 **of individuals' non-public communications)?**

86 No. For speech recordings, all the source speech signals are collected by AISHELL, which are free  
87 talk and reading. For noise recording, we are mostly recorded in public scenes and a voice activity  
88 detection method is used to filter out the recordings including prominent speech, so it avoids the  
89 contents of individuals' non-public communications.

90 **12. Does the dataset contain data that, if viewed directly, might be offensive, insulting, threaten-**  
91 **ing, or might otherwise cause anxiety?**

92 No.

93 **13. Does the dataset relate to people?**

94 No.

### 95 **B.3 Collection process**

#### 96 **1. How was the data associated with each instance acquired?**

97 The objective of the speech recording process is to mirror real-life scenarios of human activities. In  
98 each scene, the position of both the camera and microphone array are fixed. When playing source  
99 speech, the position of the loudspeaker takes on either static or moving states. For the moving case,  
100 one person manually moves the loudspeaker carrier with varying but reasonable moving speed. In  
101 transportation scenarios, people typically maintain a stationary position, thereby the loudspeaker only  
102 takes the static state. The height of the microphone array is set to 1.40 m. The center height of the  
103 loudspeaker is aligned with the height of the mouth of a standing person, varying randomly between  
104 1.30 m and 1.60 m. Most of the time, the loudspeaker faces towards the microphone array. We  
105 ensure that most speech recordings were conducted under quiet conditions (usually at midnight), with  
106 background noise levels maintained below 40 dB. Noise recording is simpler, for which we place the  
107 microphone array in various environments to capture the real-world ambient noise. Noise recording  
108 is normally conducted in the daytime with active events in each environment. Simultaneously, we  
109 have developed algorithms to perform target annotations for tasks related to speech enhancement and  
110 sound source localization.

111 **2. What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or**  
112 **sensor, manual human curation, software program, software API)?**

113 A 32-channel microphone array has been proposed for audio signals recording (including noise and  
114 speech signals played through speakers). The fisheye camera captures a 360-degree panoramic image  
115 in real time, synchronized with the microphone recording. We utilize this 360-degree panoramic  
116 image for azimuth annotation. The specific introduction of the equipment is as follows:

- 117 • 32-channel microphone array is comprised of 32 high-fidelity Audio-Technica BP899  
118 microphones. The array geometry is shown in the RealMAN paper. The audio signals  
119 are then digitized by 4 clock-synchronized 8-channel microphone pre-amplifiers (RME  
120 OctoMic II) and processed by a laptop through an audio interface (Digiface USB).
- 121 • The 360-degree fisheye camera (HIKVISION DS-2CD63C5F-IHV) is placed right above  
122 the microphone array. The frame rate of the fisheye camera is 100 ms.
- 123 • A high-fidelity monophonic loudspeaker (FOSTEX 6301 NE) is used to play source speech  
124 signals. It is placed on a height-adjustable and mobile carrier such that one can control the

125 position of the loudspeaker to mimic a standing/moving human speaker. A 5-cm diameter  
126 LED light is put on the top of the loudspeaker to magnify the visibility of loudspeaker to the  
127 the fisheye camera and annotate the position of the loudspeaker. The LED light can emit red  
128 or green light, which is visible for the fisheye camera under various of light conditions.

129 **3. If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,  
130 probabilistic with specific sampling probabilities)?**

131 N/A

132 **4. Who was involved in the data collection process (e.g., students, crowdworkers, contractors)  
133 and how were they compensated (e.g., how much were crowdworkers paid)?**

134 We did not employ external crowdworkers or contractors for data collection.

135 **5. Over what timeframe was the data collected?**

136 Data collection starts in March 2022, and ends in May 2024.

137 **6. Were any ethical review processes conducted (e.g., by an institutional review board)?**

138 No.

#### 139 **B.4 Preprocessing/cleaning/labeling**

140 **1. Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,  
141 tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing  
142 of missing values)?**

143 Yes. We provide target annotations for speech enhancement and sound source localization tasks,  
144 while also removing clearly identifiable speech segments from the noise signals.

145 **2. Was the "raw" data saved in addition to the preprocessed/cleaned/labeled data (e.g., to  
146 support unanticipated future uses)? If so, please provide a link or other access point to the  
147 "raw" data.**

148 No.

149 **3. Is the software used to preprocess/clean/label the instances available? If so, please provide a  
150 link or other access point.**

151 No.

#### 152 **B.5 Use**

153 **1. Has the dataset been used for any tasks already?**

154 Yes. In this paper, we have used it in speech enhancement (denoise and dereverberation) and speaker  
155 localization.

156 **2. Is there a repository that links to any or all papers or systems that use the dataset?**

157 Not at the present time.

158 **3. What (other) tasks could the dataset be used for?**

159 The RealMAN dataset can be used for speech enhancement and speaker localization.

160 **4. Is there anything about the composition of the dataset or the way it was collected and  
161 preprocessed/cleaned/labeled that might impact future uses?**

162 No.

163 **5. Are there tasks for which the dataset should not be used?**

164 No.

165 **B.6 Distribution**

166 **1. Will the dataset be distributed to third parties outside of the entity (e.g., company, institution,**  
167 **organization) on behalf of which the dataset was created?**

168 We have distributed the RealMAN dataset to the AISHELL website. Maybe we will distribute it to  
169 third parties.

170 **2. How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the**  
171 **dataset have a digital object identifier (DOI)?**

172 The full dataset is already publicly accessible on <https://www.aishelltech.com/RealMAN>.

173 **3. When will the dataset be distributed?**

174 It is already distributed.

175 **4. Will the dataset be distributed under a copyright or other intellectual property (IP) license,**  
176 **and/or under applicable terms of use (ToU)?**

177 The RealMAN dataset uses the CC BY 4.0 license.

178 **5. Have any third parties imposed IP-based or other restrictions on the data associated with the**  
179 **instances?**

180 No.

181 **6. Do any export controls or other regulatory restrictions apply to the dataset or to individual**  
182 **instances?**

183 No.

184 **B.7 Maintenance**

185 **1. Who is supporting/hosting/maintaining the dataset?**

186 RealMAN is hosted on AISHELL website and maintained by the AudioLab of Westlake University.

187 **2. How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

188 Contact us by email or Github.

189 **3. Is there an erratum?**

190 Not at the present time. Future errata will be published on the RealMAN Github Page.

191 **4. Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete**  
192 **instances)?**

193 We may update the dataset by adding speech recorded in more scenes to make it more valuable. Any  
194 updates will be posted updates on the RealMAN Github page.

195 **5. If the dataset relates to people, are there applicable limits on the retention of the data**  
196 **associated with the instances (e.g., were individuals in question told that their data would be**  
197 **retained for a fixed period of time and then deleted)**

198 No.

199 **6. Will older versions of the dataset continue to be supported/hosted/maintained?**

200 We expect any future changes to be additive. If that changes, we will release our versioning policy on  
201 the RealMAN Github page.

202 **7. If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for**  
203 **them to do so?**

204 We accept feedback in the issues of the RealMAN Github page.