
Supplementary Material of AV-NeRF

Anonymous Author(s)

Affiliation

Address

email

1 A Demo Videos

2 For convincing results, we have created a webpage that contains the rendering videos of multiple
3 scenes. **We highly recommend readers visit this webpage** (the “index.html” file in the folder).
4 Please note that the webpage may not be fully supported on Safari browser; therefore, we recommend
5 using Google Chrome for optimal viewing experience.

6 B Architectures

7 **A-NeRF.** A-NeRF consists of two Multilayer Perceptrons (MLPs), each comprising four linear layers
8 with an additional residual connection. The width of each linear layer, denoted as c , is set to 128 for
9 the RWAVS dataset and 256 for the SoundSpaces dataset. In A-NeRF, all linear layers are followed
10 by ReLU activation layers, except for the last layer, where the ReLU activation is replaced with the
11 Sigmoid function. The first MLP takes the listener’s position (x, y) and the frequency $f \in [0, F]$ as
12 input, where F represents the number of frequency bins. It predicts a mixture mask $m_m \in \mathcal{R}$ for
13 the given frequency f and generates a feature vector with c channels. Prior to feeding them into the
14 MLP, we apply positional encoding to the listener’s position (x, y) and the frequency f . We set the
15 maximum frequency used for positional encoding as 10.

16 Then, we adopt relative transformation (Sec. 4.4) to project the listener’s direction θ into a high-
17 frequency space. We concatenate the transformed listener’s direction and the feature vector, and
18 feed it into the second MLP. The second MLP is appended with a Sigmoid layer and a scaling layer,
19 ensuring that the difference mask m_d estimated by the second MLP falls within the range of $[-1, 1]$.
20 For each frequency query f , A-NeRF estimates two masks: m_m and m_d , both of which are scalars.
21 We iterate over all frequencies $f \in [0, F]$ to obtain the complete masks \mathbf{m}_m and \mathbf{m}_d . After computing
22 the masks, we synthesize the target audio a_t according to the procedure discussed in Sec. 4.2.

23 **V-NeRF.** We utilize the nerfacto model provided by nerf-studio [1] as the V-NeRF. This model
24 combines several well-established and successful methods, including camera pose refinement [2, 3],
25 image appearance conditioning [4], hash encoding, and proposal sampling [5]. Due to its robust and
26 effective performance on real-world data, we utilize the default settings of the nerfacto model without
27 making any architectural modifications. For more detailed information regarding the architecture of
28 V-NeRF, please refer to the documentation provided by nerf-studio.

29 **AV-Mapper.** For each camera pose, we render both RGB and depth images using V-NeRF. We resize
30 images to 256×256 and center-crop to 224×224 prior to feeding them into a frozen ResNet-18
31 [6] image encoder pre-trained on ImageNet-1K dataset [7]. ResNet-18 embeds the input image as a
32 512-dimension feature vector. We concatenate the RGB and the depth feature vectors, and input them
33 into the AV-Mapper to learn environmental knowledge of the sound acoustics. AV-Mapper projects
34 the input feature vectors to a latent embedding of c channels. We implement the AV-Mapper as a
35 3-layer MLP, with each intermediate linear layer followed by a ReLU activation function.

36 C Implementation Details

37 **RWAVS Dataset.** We implement our method using the PyTorch framework [8]. We employ Adam
38 optimizer [9] with $\beta_1 = 0.9$ and $\beta_2 = 0.999$ for model optimization. The initial learning rate is set
39 to $5e-4$ and exponentially decreased to $5e-6$. We train the model for 100 epochs with a batch size
40 of 32.

41 Before feeding the camera position (x, y) to A-NeRF, we normalize it within the range $[-1, 1] \times$
42 $[-1, 1]$ and apply positional encoding [10]. Additionally, we resample all audios to a frequency of
43 22050 Hz and utilize the Short-Time Fourier Transform (STFT) to convert waveform audios into the
44 time-frequency domain. For this transformation, we set the number of ffts as 512, the window length
45 as 512, and the hop length as 128. A Hanning window is applied during the process. Finally, we
46 compute both the magnitude and the phase from the spectrogram.

47 We use magnitude distance (MAG) [11] and envelope distance (ENV) [12] as evaluation metrics for
48 audio quality. The MAG metric quantifies audio quality in the time-frequency domain and is defined
49 as follows:

$$\text{MAG}(\mathbf{m}_{\text{prd}}, \mathbf{m}_{\text{gt}}) = \|\mathbf{m}_{\text{prd}} - \mathbf{m}_{\text{gt}}\|^2, \quad (1)$$

50 where \mathbf{m}_{prd} is the predicted magnitude, and \mathbf{m}_{gt} is the ground-truth magnitude. ENV metric that
51 measures the audio quality in the time domain is formatted as:

$$\text{ENV}(a_{\text{prd}}, a_{\text{gt}}) = \|\text{hilbert}(a_{\text{prd}}) - \text{hilbert}(a_{\text{gt}})\|^2, \quad (2)$$

52 where a_{prd} is the predicted audio, a_{gt} is the ground-truth audio, and hilbert is the Hilbert transfor-
53 mation function [13].

54 **SoundSpaces Dataset.** Our model is trained on the SoundSpaces dataset using the same training
55 settings as RWAVS dataset. We resample the impulse responses to 22050 Hz following INRAS [14].
56 The 2D position is normalized to $[-1, 1] \times [-1, 1]$ prior to positional encoding.

57 We tailor A-NeRF for impulse response prediction with some minor modifications: (1) The input
58 frequency query f is replaced by a time query $t \in [0, T]$, where T represents the length of an impulse
59 response signal; (2) the first MLP only generates a feature vector while discarding the mixture mask
60 \mathbf{m}_m ; (3) the second MLP predicts impulse response signals instead of difference mask \mathbf{m}_d .

61 Since the generated impulse responses are in the time domain, we employ STFT to convert them into
62 the time-frequency domain and calculate their magnitudes. We utilize an STFT configuration with
63 512 FFTs, a sliding window width of 512, a hop stride of 128, and a Hanning window. We supervise
64 the model training using the L2 distance between the ground-truth magnitudes and the predicted
65 magnitudes.

66 For performance evaluation, we choose three metrics: T60, C50, and EDT [14]. T60 characterizes
67 the reverberation effects in an audio signal by measuring the time it takes for the audio’s energy to
68 attenuate by 60 dB. The T60 distance is calculated as follows:

$$\text{T60}(a_{\text{prd}}, a_{\text{gt}}) = \frac{|\text{T60}(a_{\text{prd}}) - \text{T60}(a_{\text{gt}})|}{\text{T60}(a_{\text{gt}})}, \quad (3)$$

69 where a_{prd} and a_{gt} are the predicted and ground-truth impulse responses, respectively. C50 quantifies
70 the energy ratio between early reflections and late reverberation, allowing it to represent the clarity
71 and loudness of the audio. We format the C50 distance as:

$$\text{C50}(a_{\text{prd}}, a_{\text{gt}}) = |\text{C50}(a_{\text{prd}}) - \text{C50}(a_{\text{gt}})|. \quad (4)$$

72 The EDT metric shares similarities with T60 but places greater emphasis on capturing the early
73 reflections of impulse responses. The EDT distance is defined as follows:

$$\text{EDT}(a_{\text{prd}}, a_{\text{gt}}) = |\text{EDT}(a_{\text{prd}}) - \text{EDT}(a_{\text{gt}})|. \quad (5)$$

74 With these three metrics, we can evaluate the generation quality of impulse responses from different
75 aspects, including clarity, energy, and reverberation.

76 D Setup of RWAVS Dataset

77 **Recording Devices.** We have assembled a recording system, as depicted in Fig. 7, to cap-
78 ture high-quality audio-visual scenes in real-world environments. Our system comprises a
79 3Dio Free Space XLR binaural microphone for capturing stereo audio, a TASCAM DR-
80 60DMKII for recording and storing audio, and a GoPro Max for capturing accompanying videos.

81 This system is portable, allowing us to position it flexibly and capture scenes from different
82 camera poses. In addition, we utilized an LG
83 XBOOM 360 omnidirectional speaker to serve
84 as a sound source, which plays music repeatedly.
85 Figure 5 in the main paper illustrates the setup
86 used to record data in four distinct environments:
87 office, house, apartment, and outdoors. Within
88 each environment, we positioned the speaker at
89 multiple locations to capture diverse acoustic
90 effects. Each combination of environment and
91 sound source represents an audio-visual scene.
92 We collected data ranging from 10 to 25 minutes
93 for each scene, resulting in a total collection
94 of 232 minutes (3.8 hours) of diverse data, encompassing
95 various environments and source positions.
96
97



Figure 7: Recording system. It comprises a professional binaural microphone, a sports camera, and a recorder.

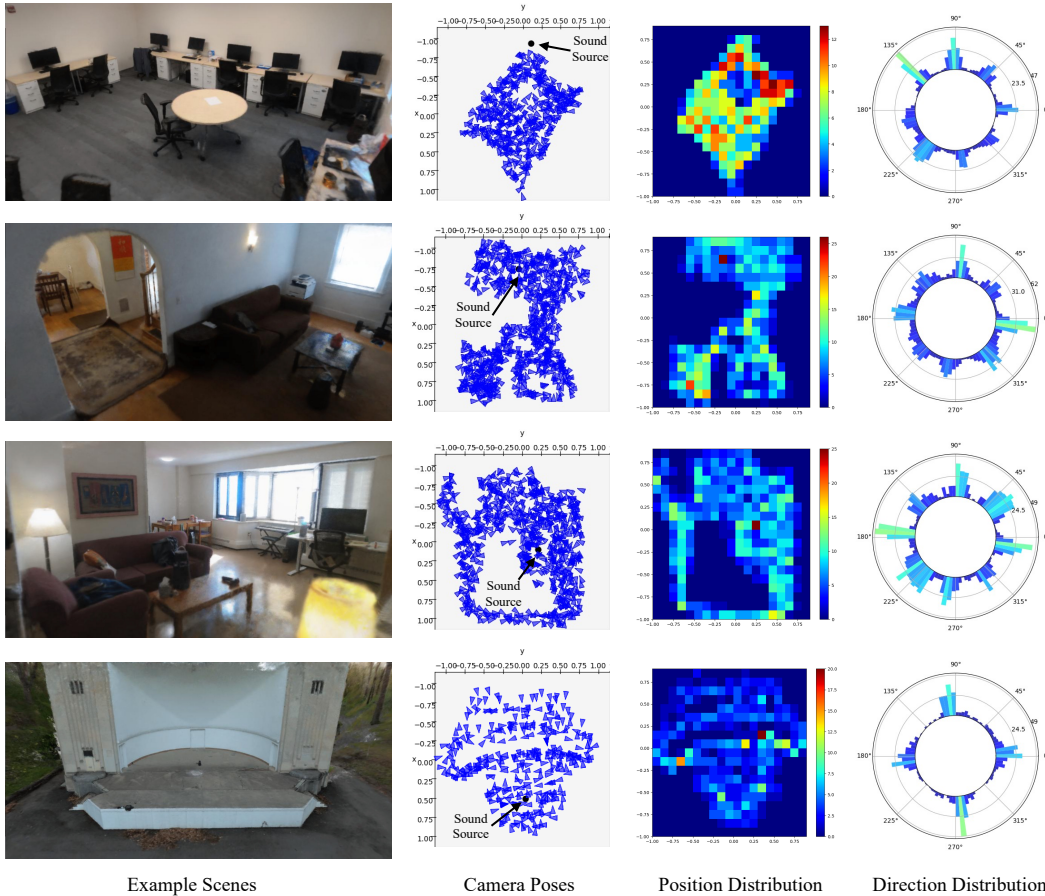


Figure 8: Example scenes. We present several example scenes along with their corresponding camera pose distributions. We display the position density heatmap and the direction distribution map. RWAVS dataset is composed of diverse environments with various camera poses.

98 **Example Scenes.** In Fig.8, we present four example scenes from RWAVS dataset along with the
99 corresponding camera pose distributions. The first column showcases images of the example scenes.
100 The second column displays the camera poses used for video recording: the black dot represents the
101 sound source, and each blue triangle represents a camera pose. We normalize all camera poses to the
102 range of $[-1, 1]$ and visualize them in an x-y plane from a top-down view. The third column contains

2D density heatmaps, which illustrate the distribution of camera poses in each unit area: each pixel represents a unit area and its color shows the number of camera poses in this area. As shown in the figure, RWAVS dataset encompasses densely covered camera poses for each environment. We also analyze the distribution of camera directions (shown in the last column of Fig.8). We present the direction distribution in a polar coordinate system with the angle representing the viewing direction and the radius meaning the number of camera poses in this region. RWAVS dataset consists of various viewpoints that approximately cover a 360° range of viewing directions. In summary, the RWAVS dataset comprises diverse environments with a wide range of camera poses.

References

- [1] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Justin Kerr, Terrance Wang, Alexander Kristoffersen, Jake Austin, Kamyar Salahi, Abhik Ahuja, David McAllister, and Angjoo Kanazawa. Nerfstudio: A modular framework for neural radiance field development. *arXiv preprint arXiv:2302.04264*, 2023.
- [2] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. NeRF—: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021.
- [3] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *ICCV*, 2021.
- [4] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. In *CVPR*, 2021.
- [5] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM TOG*, 41(4):1–15, 2022.
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009.
- [8] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *NeurIPS*, 32, 2019.
- [9] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *ICLR*, 2015.
- [10] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- [11] Xudong Xu, Hang Zhou, Ziwei Liu, Bo Dai, Xiaogang Wang, and Dahua Lin. Visually informed binaural audio generation without binaural audios. In *CVPR*, 2021.
- [12] Pedro Morgado, Nuno Vasconcelos, Timothy R. Langlois, and Oliver Wang. Self-supervised generation of spatial audio for 360° video. In *NeurIPS*, pages 360–370, 2018.
- [13] Julius Orion Smith. *Mathematics of the discrete Fourier transform (DFT): with audio applications*. 2008.
- [14] Kun Su, Mingfei Chen, and Eli Shlizerman. Inras: Implicit neural representation for audio scenes. In *NeurIPS*, 2022.