
NeRAF: 3D Scene Infused Neural Radiance and Acoustic Fields.

Amandine Brunetto Sascha Hornauer Fabien Moutarde
Centre for Robotics
Mines Paris - PSL Research University

Abstract

Sound plays a major role in human perception, providing essential scene information alongside vision for understanding our environment. Despite progress in neural implicit representations, learning acoustics that match a visual scene is still challenging. We propose NeRAF, a method that jointly learns acoustic and radiance fields. NeRAF is designed as a Nerfstudio module for convenient access to realistic audio-visual generation. It synthesizes both novel views and spatialized audio at new positions, leveraging radiance field capabilities to condition the acoustic field with 3D scene information. At inference, each modality can be rendered independently and at spatially separated positions, providing greater versatility. We demonstrate the advantages of our method on the SoundSpaces dataset. NeRAF achieves substantial performance improvements over previous works while being more data-efficient. Furthermore, NeRAF enhances novel view synthesis of complex scenes trained with sparse data through cross-modal learning. Project page: <https://amandinebttto.github.io/NeRAF>

1 Introduction

Sound is fundamental to human perception. Picture yourself in the heart of a bustling city center. It is not just what you see, but also the cacophony of car horns, footsteps, and chatter that guide your perception and decision-making. Beyond providing information about our environment, sound offers major insights into context and atmosphere, enriching our understanding of the world in ways that sight alone cannot reach. In gaming and AR/VR simulations, sound plays a crucial role in someone’s immersion, to match the environment and make things feel real.

Recent advances in novel view synthesis allowed the generation of high-quality photorealistic images at any camera position from a set of captured images [34, 33, 62, 53]. They opened many exciting applications for simulated environments. Yet these methods lack acoustic synthesis essential to immersive experiences. Learning scene acoustics presents additional challenges [30]. Indeed, sound is influenced by the geometry of the space and the materials of objects and surfaces. As sound travels from its source to our ears, it undergoes multiple acoustic phenomena. Room Impulse Responses (RIR) capture them all relative to pairs of listener and sound source positions. However, acquiring RIRs is a challenging and laborious task that involves to play and record sounds from multiple source-listener positions. Consequently, recent research [31, 47, 48] focuses on estimating RIRs at novel poses from sparsely collected data. Drawing inspiration from NeRF success in image synthesis, emerging works [30, 51, 28] learn acoustic properties of the environment using neural fields, enabling the synthesis of sound from new source and microphone locations. Still, these methods do not sufficiently account for the scene geometry and materials, despite their major role in acoustics.

In response, we introduce NeRAF, a method that synthesizes high-fidelity spatialized audios and enhances NeRF novel-view synthesis (Figure 1). NeRAF is designed as a Nerfstudio module [53] that jointly learns radiance and acoustic fields. We leverage the radiance field to obtain a 3D voxel-grid



Figure 1: NeRAF learns radiance and acoustic fields from a set of images and audio recordings. It synthesizes binaural RIRs and RGB images at novel camera, microphone and source positions and orientations. NeRAF benefits from cross-modal training without requiring co-located audio-visual data and can render spatially separate modalities. NeRAF enables auralization and audio spatialization, along with enhanced image rendering, which are essential for creating a realistic perception of space.

containing scene color and density. We encode it to condition our neural acoustic field with implicit scene geometry and materials. Similar to previous works [30, 51, 28], we validate the effectiveness of our method on Sound Spaces dataset [6, 9]. Our contributions can be summarized as follow:

- NeRAF improves binaural RIR predictions on the SoundSpaces dataset by 17.4% T60, 31.3% EDT and 31.6% C50, enabling more realistic perceptual experiences.
- Our method is data-efficient, outperforming the current state-of-the-art with only half the audio recordings.
- NeRAF supports spatially separate renderings of images and audios.
- We alleviate the constraints on the training data by allowing the use of non-co-located microphones and cameras.
- NeRAF enhances novel-view synthesis in complex scenes with sparse training data through joint learning of acoustic and radiance fields.
- We propose a novel loss combination variant that improves significantly EDT and C50.
- We provide convenient and easy access to audio-visual novel view synthesis by releasing a Nerfstudio-compatible module.

2 Related works

Neural Radiance Fields. NeRF [34] synthesizes novel photorealistic views of a scene by modeling it as a continuous function that maps 3D spatial coordinates and viewing directions to radiance. Although it has demonstrated impressive results on objects and small bounded regions, it struggles with complex scenes where the camera can be oriented in any directions and content may exist at any distance. Building upon NeRF, [33] handles unconstrained photo collections, enabling the use of diverse image sets. Other works improved NeRF to learn more complex scenes [2, 62] with higher efficiency and quality [35, 58]. Recently, Nerfstudio [53] provided a modular framework that allows for a simplified end-to-end process of creating, training, and testing NeRF. It combines many existing NeRF improvements to create Nerfacto, a model suited for real-data captures of static scenes. Our method is designed to be a Nerfstudio module bringing novel audio generation to existing state-of-the-art NeRF methods. Moreover, our method improves NeRF performance on large complex scenes given a small image set by leveraging the simultaneous neural acoustic fields training.

Neural Acoustic Fields. Several works have extended neural fields applicability to the audio domain. Luo et al. [30] (NAF) and Su et al. [51] (INRAS) model acoustic propagation in a scene by learning an implicit representation that maps emitter-listener location pairs to binaural RIRs. While NAF conditions the acoustic field by learning local geometric information, INRAS performs

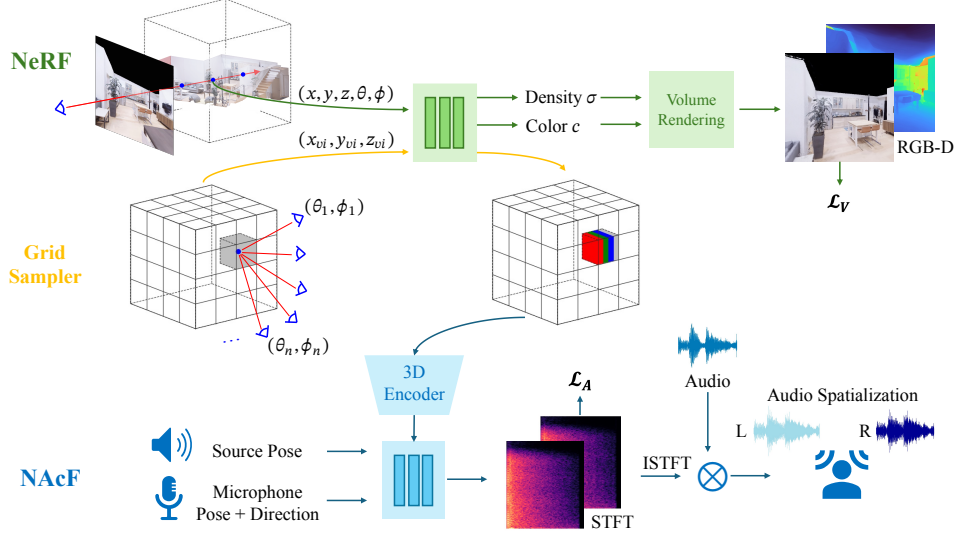


Figure 2: **Pipeline Overview.** NeRF maps 3D coordinates and orientations to density and color. The grid sampler fills a 3D grid representing the scene by querying the radiance field with voxel center coordinates and multiple viewing directions. NAcF learns to map source-microphone poses and directions to STFT coefficients. It is conditioned by extracted scene features. Predicted RIRs can be convolved with audio to obtain auralized and spatialized sound matching the scene.

audio scene feature decomposition based on the acoustic radiance transfer model. AV-NeRF [28] provides vision-based information to the acoustic field. They rely on pre-trained Nerfacto to render RGB-D images at the microphone pose and extract their features to condition binaural rendering at the same pose. Similarly, our method performs binaural RIR prediction by continuously mapping emitter-listener pairs to the acoustic field. Unlike AV-NeRF, we can decouple the camera and microphone poses used for training. This allows to benefit from unequal amounts of data sources, which is beneficial outside of a simulator. Additionally, at inference, each modality can be rendered independently for greater versatility. We go beyond camera field-of-view information by using NeRF to obtain a 3D grid of the scene filled with color and density information, bringing richer insights to the acoustic field. Furthermore, NeRF’s cross-modal training benefits both modalities, improving NeRF performances on complex scenes.

RIR synthesis. RIR have various applications, such as enabling immersive and spatialized audio for AR and VR, performing de-reverberation or providing insights about room acoustics. However, RIR collection in real-world environment is time consuming and needs specialized hardware. Consequently, synthesizing them has been a longstanding research topic. Traditional methods rely on simulated approaches such as wave-based [24, 55, 41, 3] or geometric methods [44, 46, 27, 56, 5, 49]. While these methods effectively simulate sound propagation in space, wave based methods face computation complexity and geometric methods struggle to simulate low-frequency acoustic phenomena such as interference and diffraction. Recent works have proposed learning based methods leveraging the modeling ability of neural networks to learn room acoustics and estimate RIR [31, 30, 51, 42, 43].

Audio-Visual learning. Some works rely on both audio and visual information to perform acoustic related tasks such as audio binauralization [20, 22], auralization [8, 48], de-reverberation [10, 14] or RIR prediction [47, 31]. Audio-visual learning has also demonstrated promising results for other tasks such as improving depth prediction [15, 4, 36, 63], performing floorplan reconstruction [32, 40], navigating through a space [61, 21, 13, 7, 19] and estimating poses [13].

3 Task

In this study, we focus on simulating acoustic phenomena within a static complex environment based on multiple camera images and audio recordings from known microphone and source poses

and orientations. Our goal is to simulate the audio-visual experience for novel microphone-source and camera positions. At inference, our model enables the auditory and visual exploration of the environment from any sensor position without requiring new observations. Consequently, our approach facilitates the synthesis of realistic video paired with matching spatial audio.

4 Method

Our method is designed as a Nerfstudio module that renders RGB-D images and spatialized audio at any camera and microphone-source pose and orientation. During training, it concurrently learns neural radiance and acoustic fields. At inference, each part can be used independently.

NeRAF pipeline is shown in Figure 2. It comprises 3 modules: NeRF, grid sampler and NAcF. NeRF generates RGB-D images. The grid sampler fills a 3D grid representing the scene with color and density by querying NeRF’s radiance field with the center coordinates of each voxel and multiple viewing directions. NAcF renders binaural audio conditioned by implicit materials and geometric insights provided by the encoded 3D grid.

4.1 Learning scene acoustics using RIRs

To learn the acoustic field of an environment, our method models room impulse responses. RIRs depend on scene geometry and materials. For instance, early reflections are related to the distance from nearby obstacles while late reverberation corresponds to the size and structure of the scene. Learning scene acoustic via RIRs offers several advantages: (1) RIRs capture the acoustic characteristics of a room between emitter-listener positions; (2) by convolving given sounds with RIRs, we can auralize them, effectively synthesizing sounds that mimic the perception of being in the scene; (3) binaural RIRs contain the head-related function transfer (HRTF) necessary to spatialize sound. Thus, by learning to model RIRs, our method captures the sound propagation phenomena necessary for audio synthesis with rich acoustic properties.

4.2 Neural Radiance Field

Neural radiance field (NeRF) [34] renders photo-realistic images from new view points. It represents a static scene as a 5D vector-valued continuous function whose input is a 3D location $\mathbf{X} = (x, y, z)$ and 2D viewing direction $\mathbf{d} = (\theta, \phi)$, and output is color $\mathbf{c} = (r, g, b)$ and volume density σ :

$$\text{NeRF} : (\mathbf{X}, \mathbf{d}) \rightarrow (\mathbf{c}, \sigma). \quad (1)$$

NeRF shoots rays $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ through image pixels from camera origin \mathbf{o} between the near t_{near} and far t_{far} boundaries of the scene. 3D points $\mathbf{X}_n = \mathbf{o} + t_n\mathbf{d}$ are sampled along the ray and, along with \mathbf{d} used as input of two MLPs. The first one maps \mathbf{X}_n to the view-independent density σ_n and a corresponding feature vector. The second MLP takes the feature vector and \mathbf{d} to produce view-dependent color \mathbf{c}_n . The final color $C(\mathbf{r})$ and depth $D(\mathbf{r})$ of an image pixel are rendered using volume rendering equations:

$$C(\mathbf{r}) = \int_{t_{near}}^{t_{far}} \mathcal{T}(t_n) \sigma_n \mathbf{c}_n dt_n, \quad (2)$$

$$D(\mathbf{r}) = \int_{t_{near}}^{t_{far}} \mathcal{T}(t_n) \sigma_n t_n dt_n, \quad (3)$$

where \mathcal{T} is the transmittance expressed as $\exp(-\int_{t_{near}}^{t_n} \sigma_k dk)$.

NeRAF relies on Nerfacto [53] as its NeRF model. Nerfacto is a combination of many works following [34] including camera pose refinement [57, 29], per image appearance conditioning [33], proposal sampling [35], scene contraction and hash encoding. We chose this approach because of its strong performance on real-world data captures of static scenes. However, our method is architecture-agnostic, as it only relies on the radiance field and does not modify the NeRF.

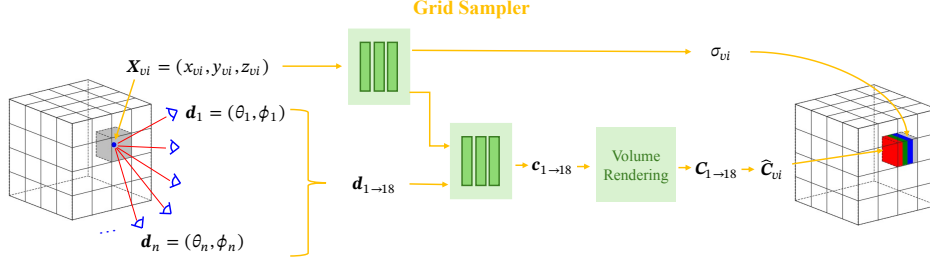


Figure 3: **Grid sampler.** We represent the scene as a grid of voxels. The grid is filled by querying the radiance field with voxel center coordinates \mathbf{X}_{vi} and multiple viewing directions $\mathbf{d}_{1 \rightarrow 18}$. We average the color values obtained for each viewing direction. It results in a 7-channels 3D grid containing color $\hat{\mathbf{C}}_{vi}$, density σ_{vi} and voxel-centers 3D coordinates.

4.3 Grid Sampler

Sound propagation is omni-directional and determined by 3D geometry and material properties of an environment, thus we provide 3D information about the scene to our neural acoustic field.

To this end, we use our grid sampler, shown in Figure 3, to construct a 3D volume from NeRF. Specifically, we represent the scene by a voxel grid of size $[0, 1]^3$ with a $1/gap$ resolution. We apply the same scene contraction as NeRF.

We obtain the positions of voxel centers \mathbf{X}_{vi} and update the grid at each training iteration using a batch of 4,096 \mathbf{X}_{vi} . This sampling approach ensures that the grid is fully populated, unlike directly using the \mathbf{X}_n sampled by NeRF proposal sampling. Indeed, NeRF learns to sample \mathbf{X}_n likely to be of interest, which could potentially leave some regions of the grid unfilled or without further updates as the model improves.

As NeRF color is view-dependent, we follow [12] and form $N = 18$ viewing rays per \mathbf{X}_{vi} . The radiance fields return a density value σ_{vi} and a color value \mathbf{c}_{vi_j} per $j \in [0, N]$ viewing directions:

$$\text{NeRF} : (\mathbf{X}_{vi}, \mathbf{d}_{1 \rightarrow N}) \rightarrow (\mathbf{c}_{vi_{1 \rightarrow N}}, \sigma_{vi}) \quad (4)$$

For color, we use the volume rendering equation and average values over all viewing directions $\hat{\mathbf{C}}_{vi}$. For density, we compute the alpha compositing value $\alpha = 1 - \exp(-\sigma_{vi}\delta)$ where δ is a chosen small value. We fill the grid with these values, resulting in 4-channels, and concatenate the \mathbf{X}_{vi} 3D coordinates to add spatial information which is important for sound propagation.

4.4 Neural Acoustic Field

The goal of our neural acoustic field (NAcF), presented in Figure 4, is to learn a continuous neural representation of the scene’s acoustics. To this end, NAcF maps microphone coordinates $\mathbf{X}_m = (x_m, y_m, z_m)$ and orientations $\mathbf{d}_m = (\theta_m, \phi_m)$ along with source position $\mathbf{X}_s = (x_s, y_s, z_s)$ and orientation $\mathbf{d}_s = (\theta_s, \phi_s)$ to a binaural room impulse response:

$$\text{NAcF} : (\mathbf{X}_m, \mathbf{d}_m, \mathbf{X}_s, \mathbf{d}_s, t) \rightarrow \text{RIR}(f, t) \quad (5)$$

Soundspaces dataset lacks the complete acoustic field parameterization described in Equation (5). Microphone rotates only around the up-axis, θ , and the source is omni-directional. Thus, NAcF is parametrized only by $(\mathbf{X}_m, \mathbf{d}_m = \theta, \mathbf{X}_s, t)$.

We encode audio using short-time Fourier transform (STFT) resulting in $\text{STFT} \in \mathbb{R}^{C \times F \times T}$ with C , F and T respectively the channels, frequency and time bins. Compared to time domain, the time-frequency representation is smoother and more suitable to neural network prediction [30]. NAcF learns the log-magnitude of the STFT. We query time t because it is particularly important for RIR: it’s crucial to perceptual effects such as reverberation. NAcF returns the frequency values corresponding to the time query. We encode t between $[-1, 1]$ and use positional encoding. Similar to NeRF [34] pose processing, microphone and source coordinates as well as direction are also embedded to a higher dimensional space using positional encoding.

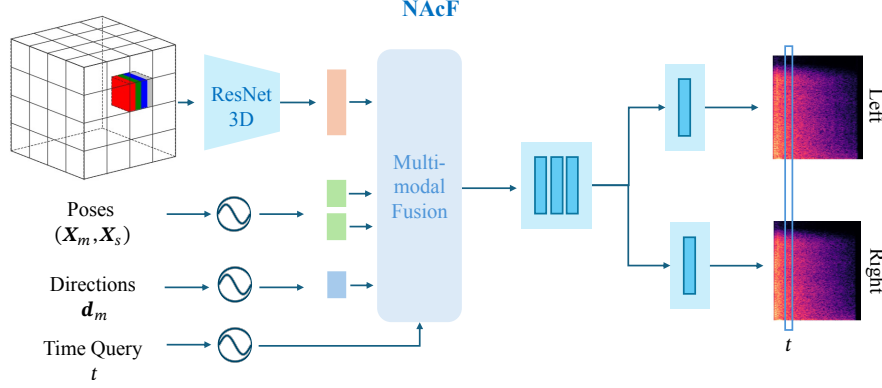


Figure 4: **Neural Acoustic Field**. NAcF maps microphone-source poses and directions to binaural RIRs. It is conditioned by scene features obtained via a 3D ResNet. Poses, directions and time queries are projected to a higher dimensional space through positional encoding. NAcF outputs a vector containing F frequency values for each time query.

Given the voxel grid, we concurrently train a ResNet3D [25] to extract pertinent features for sound propagation. They serve as implicit indicators of scene geometry and material properties which are central for acoustic modeling. This choice follows works [11, 47, 28] in audio-visual research showing that extracted geometry and material information using ResNet networks contribute to audio synthesis. We perform multimodal fusion by concatenating grid features vector with the encoded microphone and source poses, microphone direction and time query. Similar to NeRF, NAcF consists of two MLP blocks. The first block maps the input to a feature vector embedding the input information and the acoustic scene. The second learns to spatialize sound, i.e. the head related transfer function (HRTF) in the case of binaural microphones. It contains one MLP per microphone channel and output the $F = (N_{fft}/2) + 1$ frequencies for each time query.

Contrary to previous works, NAcF does not directly predict the complete STFT during training. Instead, we adopt a process similar to NeRF’s. NeRF shuffles all image pixels present in the train set and creates batches independently of their images of origin. In NeRAF, we shuffle all STFT time bins in the train set and form batches, independently of their STFT of origin. This method efficiently trains the network to learn how sound propagates over time. At inference, we query all time bins $t \in [0, T]$ to render the complete STFT.

Finally, we use Griffin-Lim [23, 38] algorithm to obtain RIR waveforms from magnitude STFTs.

4.5 Learning Objective

Cross-modal learning has proven beneficial in multiple works. Thus, NeRAF trains jointly NeRF and NAcF. NAcF training is slightly delayed as it begins when the grid has been updated a few times.

We use NeRF MSE loss function \mathcal{L}_V between ground-truth pixel color $\hat{C}(\mathbf{r})$ and the rendered one $C(\mathbf{r})$:

$$\mathcal{L}_V = \|\hat{C}(\mathbf{r}) - C(\mathbf{r})\|_2^2 \quad (6)$$

Note that each NeRF method also has complementary losses.

Inspired by [59, 60], we use a combination of spectral loss \mathcal{L}_{SL} [16] and spectral convergence loss \mathcal{L}_{SC} [1] as our audio loss \mathcal{L}_A :

$$\mathcal{L}_A = \lambda_{SC}\mathcal{L}_{SC} + \lambda_{SL}\mathcal{L}_{SL} \quad (7)$$

$$\mathcal{L}_{SL} = \|\log(|\hat{s}| + \epsilon) - \log(|s| + \epsilon)\|_2^2 \quad (8)$$

$$\mathcal{L}_{SC} = \frac{\|\hat{s} - s\|_F}{\|\hat{s}\|_F} \quad (9)$$

where \hat{s} is the ground-truth STFT, s the predicted STFT, $\|\cdot\|_F$ denotes the Frobenius norm, $\|\cdot\|_2$ the L2 norm and $\epsilon = 10^{-3}$. The spectral convergence loss emphasizes on spectral peaks helping

especially in early phases of training while the spectral loss focus on small amplitude components which tends to be more important towards the later phases of training. Our loss differs from [59, 60] as we found out that MSE spectral loss leads to better performance than L1 spectral loss. Additionally, our loss combination outperforms the MSE loss on magnitude STFT used in [30, 28]. The loss comparison is presented in Section 5.3.

NeRAF complete learning objective can be resumed as:

$$\mathcal{L} = \mathcal{L}_V + \lambda_A \mathcal{L}_A \quad (10)$$

5 Experiments

First, we evaluate NeRAF’s performance against NAF, INRAS and AV-NeRF on the SoundSpaces dataset. Then, we demonstrate that cross-modal learning of the acoustic and radiance fields is beneficial to vision, resulting in enhanced results for large complex scenes with very sparse training observations. Next, we conduct a few-shot experiment to assess our method with a reduced number of training audio recordings. Finally, we investigate the impact of the 3D voxel-grid and of our loss combination on NeRAF.

5.1 Datasets

To the best of our knowledge, there is no existing real-world dataset that includes RGB images with camera poses and orientations, along with RIRs and microphone-source poses. Especially, for NeRF to perform well, training views must sufficiently cover the scene. Therefore, we follow previous works and evaluate NeRAF on the SoundSpaces simulated dataset.

SoundSpaces. SoundSpaces [6] is an audio simulator built upon Habitat Sim [45, 52, 39], a 3D simulator for embodied AI research. It provides binaural RIRs at discrete positions of a 2D grid with a spatial resolution of 0.5 m with four head orientations accompanied by sound source position. This popular framework [21, 7, 30, 28, 51, 40, 36, 8] enables us to evaluate our method on diverse scene types while being under same settings as previous works. Following [30, 28, 51] we select six indoor scenes from Replica [50]. Two are small rooms (office 4 and room 2) with rectangular walls, two are medium room (frrl apartment 2 and 4) with more complex layout and objects and two are complete apartments containing multiple rooms (apartment 1 and 2). For each scene, we use SoundSpaces pre-generated binaural RIRs with corresponding source and microphone positions. To train NeRF, we generate RGB images using Habitat Sim. Camera poses and parameters were selected to limit the amount of visual data. More details on the motivations and on the generation process are available in the supplementary material.

5.2 Evaluation

Metrics. Following [51, 28], we assess the quality of the predicted impulse responses using Reverberation Time (T60), Clarity (C50) and Early Decay Time (EDT). T60 reflects the overall sound decay within a room by measuring the time an impulse response takes to decay by 60 dB. C50 measures the energy ratio between the first 50 ms of RIR and the remaining portion. It relates to speech intelligibility, and the clarity of acoustics. EDT focuses on early reflections.

Table 1: **Comparison with State-of-the-art.** Performance on the SoundSpaces dataset using T60, C50, and EDT metrics. Lower score indicates higher RIR quality.

Methods	T60 Error (%) ↓	C50 Error (dB) ↓	EDT Error (sec) ↓
Opus-nearest	10.10	3.58	0.115
Opus-linear	8.64	3.13	0.097
AAC-nearest	9.35	1.67	0.059
AAC-linear	7.88	1.68	0.057
NAF [30]	3.18	1.06	0.031
INRAS [51]	3.14	0.60	0.019
AV-NeRF [28]	2.47	0.57	0.016
NeRAF	2.04	0.39	0.011

Table 2: **Quantitative Results on Cross-Modal Learning.** For highly complex scenes with very sparse training images, NeRAF audio-visual joint training improves vision performances, especially LPIPS. Evaluation on 50 test novel views. Delta are computed against Nerfacto (no sound). Metrics are averaged over the two large rooms.

Training Images	PSNR \uparrow			SSIM \uparrow			LPIPS \downarrow		
	75	100	150	75	100	150	75	100	150
Δ NeRAF	+3.62%	+2.24%	+1.23%	+1.44%	+1.53%	+1.03%	-17.11%	-9.86%	-6.25%

To assess the effect of cross-modal learning on novel view synthesis, we employ commonly used metrics: PSNR, SSIM, and LPIPS.

Experimental Setup. Similar to [51, 30, 28], we use 90% of SoundSpaces audio data for training and 10% for testing. We resample the RIR to 22,050 Hz and compute STFT using $N_{fft} = 512$ and $hop = 128$. To train the vision part, we generate images using SoundSpaces 2.0 [9]. For small, medium and large rooms, we use respectively 45, 75 and 150 observations. Evaluation set comprises 50 images.

Results. We compare NeRAF with state-of-the-arts methods on SoundSpaces in Table 1. Similar to previous works on neural acoustic field, we also compare our method to traditional audio encoding methods AAC and Opus [17, 18]. AAC is a multichannel audio coding standard and Opus is an open audio codec. With a grid of resolution 128, NeRAF significantly outruns existing methods. It achieves 17.4% improvement on T60, 31.3% on EDT and 31.6% on C50 compared to AV-NeRF. Since they also employ Nerfacto, our performance improvement is not due to a superior NeRF architecture, but rather to our method’s design. We showcase qualitative results of audio synthesis in supplementary material.

Cross-Modal Learning. We next explore the effect of jointly learning acoustics and radiance fields on complex scenes with a low amount of training views. We train NeRAF and Nerfacto on the two large rooms. 75, 100, 150 images are used for training and the same 50 novel views are used for evaluation. In Table 2 we show that joint training of acoustic and radiance field improves novel view synthesis. It has a stronger impact on LPIPS which is related to the human-perceived similarity between two images. Overall, the performance improvement is more significant when the training images are very sparse. We present qualitative results in Figure 5.



Figure 5: **Visualization of Cross-Modal Learning Impact on Vision Performances.** We compare Nerfacto (Vision Only) with NeRAF (Audio + Vision) on apartment 1 and 2. Learning both radiance and acoustic fields improves the quality of novel views synthesis, reducing floating artifacts and enabling more detailed images.

Few-shot RIR Synthesis. Previous works rely on a significant amount of RIR data, which is often unavailable in real-life scenarios. Therefore, we assess the performance of our model when trained with fewer data. We conduct a few-shot experiment on one small, one medium, and one large room and average the results. As shown in Figure 6, with only half the data, NeRAF outperforms AV-NeRF trained on the complete training set.

5.3 Ablation Study

Grid Impact. We compare NeRAF performances without grid and with grids of resolution 128^3 and 256^3 . Results are presented in Table 3 (Left). Without the voxel grid, NAcF is no longer conditioned by the 3D information from the scene. It still outperforms AV-NeRF with 10.5%, 28.1% and 31.3% improvement on T60 error, C50 and EDT. Compared to no grid, the 128-resolution 3D grid improves T60 error by 7.5% and C50 by 5.3%. We do not observe an enhancement in EDT. However, increasing the resolution to 256 improves EDT by 5.9%. EDT is related to RIR early reflections which usually corresponds to obstacles very close to the listener. Thus, it makes sense that

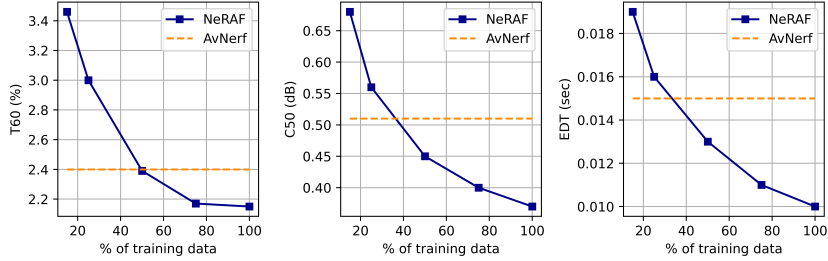


Figure 6: **Few-shot Learning.** Average on office 4, frl apartment 2 and apartment 2. With just half the data, NeRAF outperforms AV-NeRF trained on all the training data.

Table 3: **Ablation study.** Bold and underlined respectively indicates best and second best. **Left:** Influence of the grid. Evaluation on all scenes in SoundSpaces dataset. NG stands for no grid, while 128 and 256 respectively correspond to grid resolutions 128^3 and 256^3 . **Right:** Loss study. Performances averaged across office 4, frl apartment 2 and apartment 2.

Methods	T60 (%) ↓	C50 (dB) ↓	EDT (sec) ↓	Methods	T60 (%) ↓	C50 (dB) ↓	EDT (sec) ↓
NG	<u>2.21</u>	0.41	0.011 ₄	MSE	2.20	<u>0.52</u>	<u>0.014</u>
128	2.04	<u>0.39</u>	0.011 ₄	SC+SL _{L1}	1.98	0.57	0.016
256	2.38	0.38_g	0.010₇	Ours	<u>2.15</u>	0.37	0.010

the method needs greater resolution i.e. smaller voxels to perceive those obstacles. Increasing the grid resolution also improves C50, but we observe a decrease in T60 performance. T60 is related to the late reverberation of the RIR, which depends mainly on the size of the room. Consequently, increasing grid resolution should not significantly improve the quality of this information.

Loss impact. We examine the impact of our custom loss described in Section 4.5 compared to the MSE loss used in previous works [28, 30] and the combination with L1 spectral loss proposed by [60]. To this end, we evaluate NeRAF’s performance using these losses. As shown in Table 3 (Right), compared to the MSE loss, our combination significantly improves EDT and C50 by 29.1% and 27.9%, respectively, and slightly enhances T60 by 2.21%. Additionally, using MSE for \mathcal{L}_{SL} instead of L1 substantially improves C50 and EDT by 34.1% and 36.4%, although it results in 9% decline in T60.

6 Discussion

Limitation and Future Work. A primary limitation of our method is its scene dependency, requiring NeRAF to be trained separately for each scene. Future work should explore advancements in generalizable NeRF to develop a unified, generalizable approach for synthesizing audio-visual scenes. Additionally, due to the absence of a suitable real-world dataset for our task, our method has not been tested on real data. Future research will address this by evaluating the method in real-world scenarios. The current approach is also limited to static scenes with a single sound source at a time. Addressing the challenge of learning implicit representations for audio-visual scenes with multiple and dynamic sound sources would be a significant advancement. Finally, our work relies on RIRs, which, while effective at capturing scene acoustics, are not easy to obtain. It would be beneficial to evaluate the method using more common or ambient sounds.

Societal Impact. By synthesizing realistic audio-visual scenes, NeRAF can enable immersive experiences in VR and gaming. Understanding acoustics is also useful for applications such as sound de-reverberation, source localization, and agent navigation. However, if misused, our method could contribute to the creation of misleading media, including the ability to mask and lie about someone’s location.

Conclusion. We introduce NeRAF, a cross-modal method that learns both neural radiance and acoustic fields. By conditioning the acoustic field with scene features from a 3D voxel grid containing color and density information from NeRF, NeRAF enables realistic audio auralization, spatialization and novel view synthesis. Our results demonstrate its promise: substantial improvement over other methods, enhanced novel view synthesis in complex scenes, increased data efficiency and easier access to audio-visual synthesis.

Acknowledgments and Disclosure of Funding

We would like to thank Simon de Moreau for his help and insights. This work was supported by the French Agence Nationale de la Recherche (ANR), under grant ANR22-CE94-0003.

References

- [1] S. Ö. Arık, H. Jun, and G. Diamos. Fast spectrogram inversion using multi-head convolutional neural networks. *IEEE Signal Processing Letters*, 2018.
- [2] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *CVPR*, 2022.
- [3] S. Bilbao, A. Politis, and B. Hamilton. Local time-domain spherical harmonic spatial encoding for wave-based acoustic simulation. *IEEE Signal Processing Letters*, 2019.
- [4] A. Brunetto, S. Hornauer, X. Y. Stella, and F. Moutarde. The audio-visual batvision dataset for research on sight and sound. In *IROS*, 2023.
- [5] C. Cao, Z. Ren, C. Schissler, D. Manocha, and K. Zhou. Interactive sound propagation with bidirectional path tracing. *ACM TOG*, 2016.
- [6] C. Chen, U. Jain, C. Schissler, S. V. A. Gari, Z. Al-Halah, V. K. Ithapu, P. Robinson, and K. Grauman. Soundspaces: Audio-visual navigation in 3d environments. In *ECCV*, 2020.
- [7] C. Chen, S. Majumder, Z. Al-Halah, R. Gao, S. K. Ramakrishnan, and K. Grauman. Audio-visual waypoints for navigation. *CoRR*, abs/2008.09622, 2020.
- [8] C. Chen, R. Gao, P. Calamia, and K. Grauman. Visual acoustic matching. In *CVPR*, 2022.
- [9] C. Chen, C. Schissler, S. Garg, P. Kobernik, A. Clegg, P. Calamia, D. Batra, P. W. Robinson, and K. Grauman. Soundspaces 2.0: A simulation platform for visual-acoustic learning. In *NeurIPS*, 2022.
- [10] C. Chen, W. Sun, D. Harwath, and K. Grauman. Learning audio-visual dereverberation. In *ICASSP*, 2023.
- [11] C. Chen, W. Sun, D. Harwath, and K. Grauman. Learning audio-visual dereverberation. In *ICASSP*, 2023.
- [12] Y. Chen and G. H. Lee. Dreg-nerf: Deep registration for neural radiance fields. In *ICCV*, 2023.
- [13] Z. Chen, S. Qian, and A. Owens. Sound localization from motion: Jointly learning sound direction and camera rotation. In *ICCV*, 2023.
- [14] S. Chowdhury, S. Ghosh, S. Dasgupta, A. Ratnarajah, U. Tyagi, and D. Manocha. Adverb: Visually guided audio dereverberation. In *ICCV*, 2023.
- [15] J. H. Christensen, S. Hornauer, and X. Y. Stella. Batvision: Learning to see 3d spatial layout with two ears. In *ICRA*, 2020.
- [16] A. Défossez, N. Zeghidour, N. Usunier, L. Bottou, and F. Bach. Sing: Symbol-to-instrument neural generator. *NeurIPS*, 2018.
- [17] I. O. for Standardization. Advanced audio coding (aac), 2006.
- [18] X. O. Foundation. Xiph opus. <https://opus-codec.org/>, 2012.
- [19] C. Gan, Y. Zhang, J. Wu, B. Gong, and J. B. Tenenbaum. Look, listen, and act: Towards audio-visual embodied navigation. In *ICRA*, 2020.
- [20] R. Gao and K. Grauman. 2.5 d visual sound. In *CVPR*, 2019.
- [21] R. Gao, C. Chen, Z. Al-Halah, C. Schissler, and K. Grauman. Visualechoes: Spatial image representation learning through echolocation. In *ECCV*, 2020.

- [22] R. Garg, R. Gao, and K. Grauman. Visually-guided audio spatialization in video with geometry-aware multi-task learning. *IJCV*, 2023.
- [23] D. Griffin and J. Lim. Signal estimation from modified short-time fourier transform. *IEEE Transactions on acoustics, speech, and signal processing*, 1984.
- [24] N. A. Gumerov and R. Duraiswami. A broadband fast multipole accelerated boundary element method for the three dimensional helmholtz equation. *ASA*, 2009.
- [25] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [26] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [27] A. Krokstad, S. Strom, and S. Sørsdal. Calculating the acoustical room response by the use of a ray tracing technique. *JSV*, 1968.
- [28] S. Liang, C. Huang, Y. Tian, A. Kumar, and C. Xu. Av-nerf: Learning neural fields for real-world audio-visual scene synthesis. In *NeurIPS*, 2023.
- [29] C.-H. Lin, W.-C. Ma, A. Torralba, and S. Lucey. Barf: Bundle-adjusting neural radiance fields. In *ICCV*, 2021.
- [30] A. Luo, Y. Du, M. Tarr, J. Tenenbaum, A. Torralba, and C. Gan. Learning neural acoustic fields. In *NeurIPS*, 2022.
- [31] S. Majumder, C. Chen, Z. Al-Halah, and K. Grauman. Few-shot audio-visual learning of environment acoustics. *NeurIPS*, 2022.
- [32] S. Majumder, H. Jiang, P. Moulon, E. Henderson, P. Calamia, K. Grauman, and V. K. Ithapu. Chat2map: Efficient scene mapping from multi-ego conversations. In *CVPR*, 2023.
- [33] R. Martin-Brualla, N. Radwan, M. S. Sajjadi, J. T. Barron, A. Dosovitskiy, and D. Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *CVPR*, 2021.
- [34] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- [35] T. Müller, A. Evans, C. Schied, and A. Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM TOG*, 2022.
- [36] K. K. Parida, S. Srivastava, and G. Sharma. Beyond image to depth: Improving depth prediction using echoes. In *CVPR*, 2021.
- [37] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *NeurIPS*, 2019.
- [38] N. Perraudin, P. Balazs, and P. L. Søndergaard. A fast griffin-lim algorithm. In *2013 IEEE workshop on applications of signal processing to audio and acoustics*, 2013.
- [39] X. Puig, E. Undersander, A. Szot, M. D. Cote, R. Partsey, J. Yang, R. Desai, A. W. Clegg, M. Hlavac, T. Min, T. Gervet, V. Vondrus, V.-P. Berges, J. Turner, O. Maksymets, Z. Kira, M. Kalakrishnan, J. Malik, D. S. Chaplot, U. Jain, D. Batra, A. Rai, and R. Mottaghi. Habitat 3.0: A co-habitat for humans, avatars and robots, 2023.
- [40] S. Purushwalkam, S. V. A. Gari, V. K. Ithapu, C. Schissler, P. Robinson, A. Gupta, and K. Grauman. Audio-visual floorplan reconstruction. In *ICCV*, 2021.
- [41] N. Raghuvanshi, R. Narain, and M. C. Lin. Efficient and accurate sound propagation using adaptive rectangular decomposition. *TVCG*, 2009.
- [42] A. Ratnarajah, Z. Tang, and D. Manocha. IR-GAN: Room Impulse Response Generator for Far-Field Speech Recognition. In *Interspeech*, 2021.

- [43] A. Ratnarajah, S.-X. Zhang, M. Yu, Z. Tang, D. Manocha, and D. Yu. Fast-rir: Fast neural diffuse room impulse response generator. In *ICASSP*, 2022.
- [44] L. Savioja and U. P. Svensson. Overview of geometrical room acoustic modeling techniques. *ASA*, 2015.
- [45] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijmans, B. Jain, J. Straub, J. Liu, V. Koltun, J. Malik, D. Parikh, and D. Batra. Habitat: A Platform for Embodied AI Research. In *ICCV*, 2019.
- [46] C. Schissler and D. Manocha. Interactive sound propagation and rendering for large multi-source scenes. *ACM TOG*, 2016.
- [47] N. Singh, J. Mentch, J. Ng, M. Beveridge, and I. Drori. Image2reverb: Cross-modal reverb impulse response synthesis. In *ICCV*, 2021.
- [48] A. Somayazulu, C. Chen, and K. Grauman. Self-supervised visual acoustic matching. *NeurIPS*, 2024.
- [49] T. Sprunck, A. Deleforge, Y. Privat, and C. Foy. Gridless 3d recovery of image sources from room impulse responses. *IEEE Signal Processing Letters*, 2022.
- [50] J. Straub, T. Whelan, L. Ma, Y. Chen, E. Wijmans, S. Green, J. J. Engel, R. Mur-Artal, C. Ren, S. Verma, A. Clarkson, M. Yan, B. Budge, Y. Yan, X. Pan, J. Yon, Y. Zou, K. Leon, N. Carter, J. Briaies, T. Gillingham, E. Mueggler, L. Pesqueira, M. Savva, D. Batra, H. M. Strasdat, R. D. Nardi, M. Goesele, S. Lovegrove, and R. Newcombe. The Replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019.
- [51] K. Su, M. Chen, and E. Shlizerman. Inras: Implicit neural representation for audio scenes. In *NeurIPS*, 2022.
- [52] A. Szot, A. Clegg, E. Undersander, E. Wijmans, Y. Zhao, J. Turner, N. Maestre, M. Mukadam, D. Chaplot, O. Maksymets, A. Gokaslan, V. Vondrus, S. Dharur, F. Meier, W. Galuba, A. Chang, Z. Kira, V. Koltun, J. Malik, M. Savva, and D. Batra. Habitat 2.0: Training home assistants to rearrange their habitat. In *NeurIPS*, 2021.
- [53] M. Tancik, E. Weber, E. Ng, R. Li, B. Yi, J. Kerr, T. Wang, A. Kristoffersen, J. Austin, K. Salahi, A. Ahuja, D. McAllister, and A. Kanazawa. Nerfstudio: A modular framework for neural radiance field development. In *SIGGRAPH*, 2023.
- [54] D. Thery and B. F. Katz. Anechoic audio and 3d-video content database of small ensemble performances for virtual concerts. In *Intl Cong on Acoustics (ICA)*, 2019.
- [55] L. L. Thompson. A review of finite-element methods for time-harmonic acoustics. *ASA*, 2006.
- [56] M. Vorländer. Simulation of the transient and steady-state sound propagation in rooms using a new combined ray-tracing/image-source algorithm. *ASA*, 1989.
- [57] Z. Wang, S. Wu, W. Xie, M. Chen, and V. A. Prisacariu. Nerf-: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021.
- [58] Q. Xu, Z. Xu, J. Philip, S. Bi, Z. Shu, K. Sunkavalli, and U. Neumann. Point-nerf: Point-based neural radiance fields. In *CVPR*, 2022.
- [59] R. Yamamoto, E. Song, and J.-M. Kim. Probability density distillation with generative adversarial networks for high-quality parallel waveform generation. *Interspeech*, 2019.
- [60] R. Yamamoto, E. Song, and J.-M. Kim. Parallel wavegan: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram. In *ICASSP*, 2020.
- [61] A. Younes, D. Honerkamp, T. Welschehold, and A. Valada. Catch me if you hear me: Audio-visual navigation in complex unmapped environments with moving sounds. *RA-L*, 2023.
- [62] K. Zhang, G. Riegler, N. Snavely, and V. Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020.
- [63] L. Zhu, E. Rahtu, and H. Zhao. Beyond visual field of view: Perceiving 3d environment with echoes and vision. *arXiv preprint arXiv:2207.01136*, 2022.

Supplementary Material

A Architectures

NeRF. NeRAF works with any NeRF method without modifications. It only requires a radiance field that the grid sampler can query using voxel center coordinates and viewing directions. Our method only relies on NeRF to obtain a voxel-grid representation of the scene. In this paper, we choose Nerfacto from Nerfstudio [53]. For more details on its architecture we invite readers to refer to Nerfstudio website and paper.

Grid Sampler. The grid sampler creates a 3D grid of voxels with a given resolution. Grid coordinates are comprised between $[0, 1]^3$ fitting the scene contraction performed by NeRF. Consequently, knowing the exact dimensions of the room is not necessary. Voxels center coordinates and 18 viewing directions are queried to NeRF with no architecture modification. The grid is filled with color and density information along with voxels coordinates.

NACF. We train a ResNet3D-50 to embed the grid into 2,048 features that best condition the acoustic field. At inference, only these features, and not the complete grid, are required to perform RIR synthesis. Note that we employ average pooling at the end of the ResNet, adjusting according to the grid size to consistently yield 2,048 features. NACF consists of 2 MLP blocks. The first MLP block takes as input the multimodal fusion of grid features, encoded positions, and directions obtained through vector concatenation. It comprises 5 layers, each followed by a Leaky ReLU activation with a slope of 0.1, and outputs a 512 intermediate representation of the acoustic field. The second MLP block takes this intermediate representation as input and predicts the F frequencies corresponding to the time query, aiming to learn the HRTF. It consists of one separate MLP layer per microphone channel. The final activation layer is a tanh function scaled between $[-10, 10]$ to fit the STFT log-magnitude range.

B Implementation Details

We implement our method using PyTorch framework [37]. We optimize NACF using Adam optimizer [26] with $\beta_1 = 0.9$ and $\beta_2 = 0.999$ and $\epsilon = 10^{-15}$. The initial learning rate is 10^{-4} . It decreases exponentially to reach 10^{-8} . For NeRF, we keep default Nerfacto parameters.

For the first 2,000 iterations, we only train the NeRF part. It allows the grid to be filled and updated several times using batches of 4,096 voxel-centers. After, both NeRF and NACF are train jointly. We use batch sizes of 4,096 for NeRF and 2,048 for NACF. NeRAF is trained for 500k iterations but most runs reach their peak performance before, depending of the room size. We train our method on a single RTX 4090 GPU.

Similar to previous works, we resample all RIRs from 44,100 Hz to 22,050 Hz. We compute STFT with 512 FFT bins, a Hann window of size 512 and hop length of 128. We obtain log-magnitude STFT using $\log(|\text{STFT}| + 10^{-3})$. Like [30], we cut STFT to a maximum length depending of the scene. If the STFT is shorter than the maximum length we pad it with its minimal value. We inverse the transform of the predicted STFT using PyTorch Griffin-Lim algorithm running on GPU. We compute metrics against ground-truth waveforms.

We give 3D microphone and source positions to NACF. They go through multi-scale positional encoding with $N = 10$ frequencies and maximum frequency exponent of $e = 8$. Microphone rotation is scaled between $[0, 2\pi]$. Then, we apply positional encoding with $N = 4$ and $e = 8$. The time queries are normalized within the range $[-1, 1]$ and go through positional encoding with $N = 10$ and $e = 8$.

For the learning objective defined in Section 4.5, we select $\lambda_A = 10^{-3}$, $\lambda_{SC} = 10^{-1}$ and $\lambda_{SL} = 1$.

To evaluate the quality of synthesized RIR waveforms, w , we use T60, C50 and EDT errors obtained as follows:



Figure S1: Locations of visual observations for 3 SoundSpaces scenes. Blue dots corresponds to training views and green crosses to evaluation views.

$$T60(w, w_{GT}) = \frac{|T60(w) - T60(w_{GT})|}{T60(w_{GT})} \quad (11)$$

$$C50(w, w_{GT}) = |C50(w) - C50(w_{GT})| \quad (12)$$

$$EDT(w, w_{GT}) = |EDT(w) - EDT(w_{GT})| \quad (13)$$

To evaluate NeRF performances we choose metrics commonly used such as PSNR, SSIM and LPIPS.

C Scene Views Generation

SoundSpaces 1.0 comes with a code to generate RGB and depth observations at microphone positions and orientations. They corresponds to 128×128 resolution images with a 60° field-of-view (FOV) sampled every 0.5 m with orientations $\in [0^\circ, 90^\circ, 180^\circ, 270^\circ]$. In NeRAF, we decided to not rely on them for NeRF training and instead generate our own observations using Habitat Sim. This is motivated by real scenarios constraints. Views of the space may be obtained via a video recorded in the room. It is unlikely that someone will stop every 0.5 m and turn around to obtain images at each orientations. Moreover, those kind of views are not optimal for NeRF training leading, in the worst case, to the failure of less sophisticated methods, and, in the best case, to the use of more data. We also wanted to decorrelate microphone poses and camera poses for greater freedom in the usage of the method.

With Habitat Sim we rendered higher resolutions images with size 512×512 and with a 90° FOV. The larger FOV ensures more overlap between views. We randomly select positions at the edge of the room, orienting the camera to face its center with a random offset. For small room, medium and large room we use respectively 45, 75 and 150 training images. 50 test poses are randomly sampled in the room. Figure S1 presents examples of the positions obtained.

D Number of Parameters, Storage and Inference Speed

We provide in Table S1 the numbers of parameters of the NeRAF model and evaluate its inference storage requirements and speed. The speed was averaged over the 6 rooms on a single RTX 4090. The storage is the same for every rooms as they all use a 2,048-features representation of the scene regardless of its size. Note that at inference, the grid is not necessary and the extracted features are sufficient allowing the method to be faster and more compact.

Table S1: Model number of parameters along with storage and speed requirements at inference.

Parameters (Million)	Storage (MB)	Speed (ms)
71.85	97.70	7.04

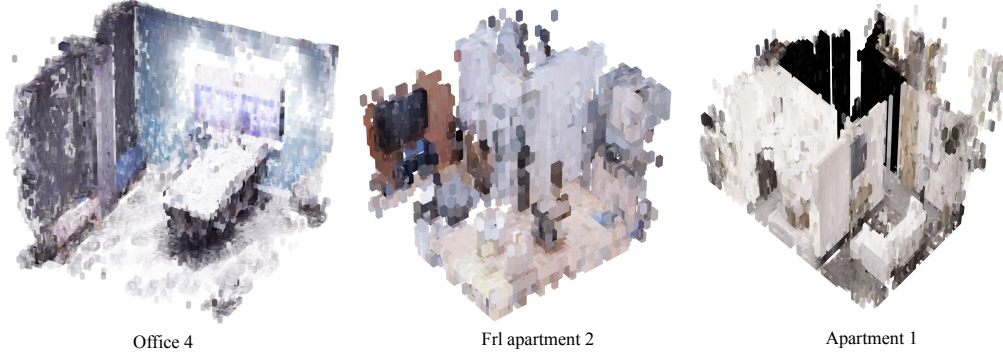


Figure S2: Sectional views of 128^3 resolution grids obtained using the grid sampler.

E Grid Visualization

We display examples of sectional views of grids obtained using the grid sampler in Figure S2. They are then encoded into features using ResNet, which condition our neural acoustic fields.

F Additional Results Visualization

Audio qualitative results. We provide examples of log-magnitude STFTs predicted with NeRAF in Figure S3

Distance-aware spatialized audio. We render binaural RIRs at various distances from the sound source and perform auralization by convolving an anechoic sound with these RIRs. This demonstrates one of the primary applications of NeRAF. The chosen sound is open source and sourced from AVAD-VR [54]. In Figure S4, we illustrate that NeRAF is distance-aware: as the microphone approaches the sound source, the amplitude increases, and it decreases as the microphone moves away. Additionally, when the microphone is positioned to the right or left of the source, the audio channels' amplitudes reflect this spatial positioning.

Cross-modal learning. We provide additional qualitative results on the effect of joint learning neural acoustic and radiance fields in Figure S5.

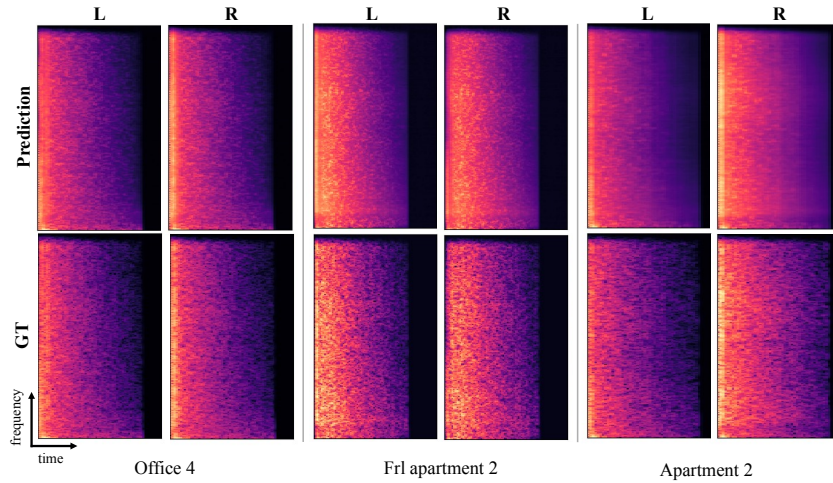


Figure S3: Examples of predicted log-magnitudes of binaural RIR STFTs compared to corresponding ground truth (GT). L and R respectively stand for left and right.

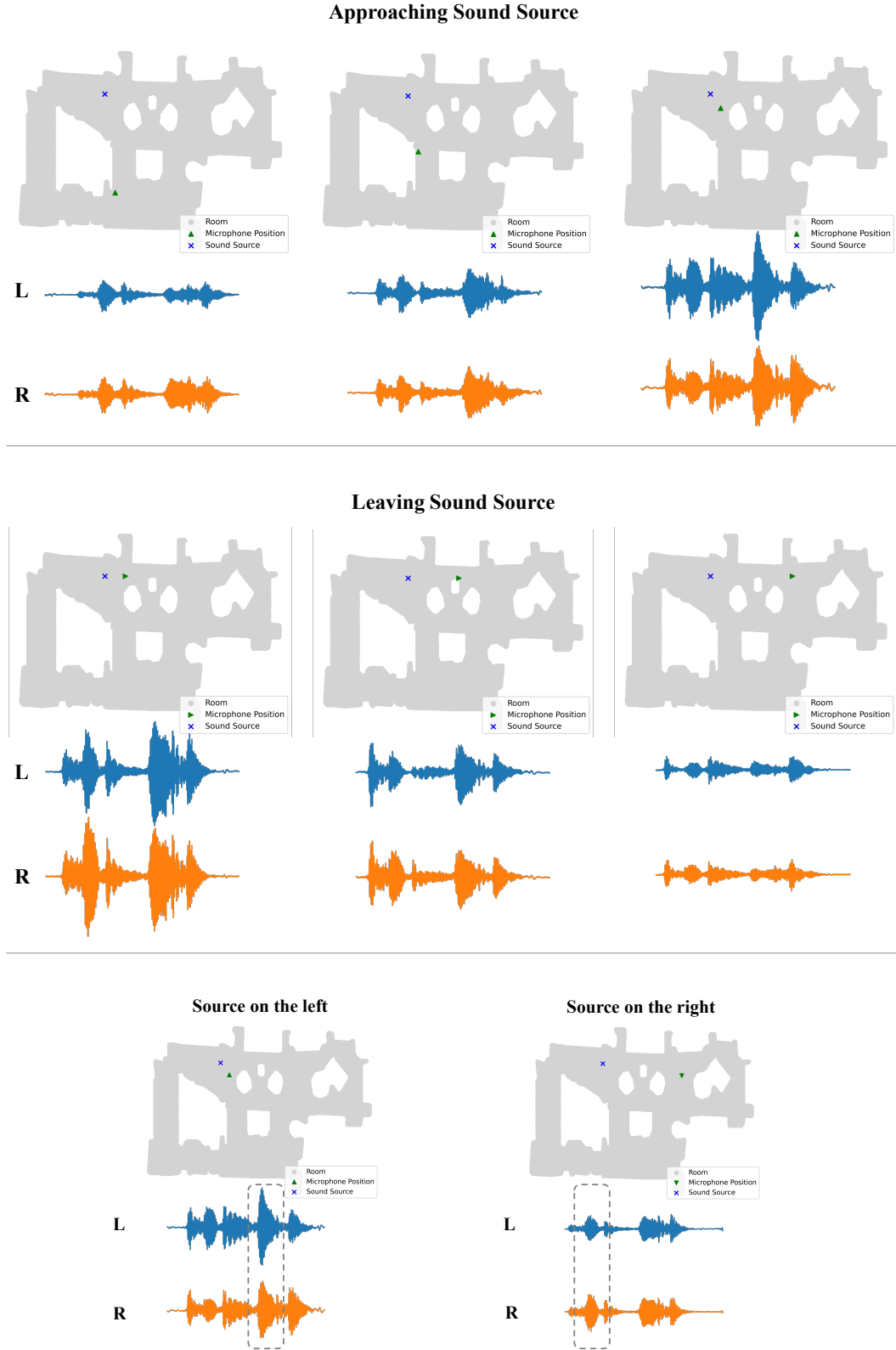


Figure S4: Examples of sounds generated with NeRAF on frl apartment 2. Our method successfully synthesizes spatialized and distance-aware sounds.



Figure S5: Cross-modal learning additional visualization.

G License Information

Asset	License
Soundspaces [6]	Creative Commons Attribution 4.0 International License
Habitat-Sim [45]	MIT License
Nerfstudio [53]	Apache License 2.0
Replica [50]	Replica license