

---

# Supplemental Material for Sounding Bodies: Modeling 3D Spatial Sound of Humans Using Body Pose and Audio

---

**Xudong Xu\***

Shanghai AI Laboratory  
xudongxu9710@gmail.com

**Dejan Marković**

Meta Reality Labs Research  
dejanmarkovic@meta.com

**Jacob Sandakly**

Meta Reality Labs Research  
jasandakly@meta.com

**Todd Keebler**

Meta Reality Labs Research  
toddkeebler@meta.com

**Steven Krenn**

Meta Reality Labs Research  
stevenkrenn@meta.com

**Alexander Richard**

Meta Reality Labs Research  
richardalex@meta.com

## A Time Warping

As described in Section 4.2, we employ time warping to tackle the time delay caused by sound propagation. While this delay in principle can be learned by the network, it has been shown that neural spatial audio modeling benefits from explicit time warping (*i.e.*, modeling the delay geometrically by shifting the input signals), as without it the model capacity is spent on learning the delay, causing slow convergence and less accurate results Richard et al. [2021]. While in tasks with a single and known sound source location, geometric time warping is trivial, in our case the origin of the sound is unknown, in fact it is even possible that there are multiple sound sources like left hand and right hand snapping simultaneously. To overcome the issue, we identify six body joints (hands, feet, nose, and hip) that are the most frequent sound emitters, and warp the sound from each of these positions towards the target microphones.

The pipeline of time warping is illustrated in Figure 1. For simplicity, we describe the warping process for a single input microphone (1-channel input). In practice, we apply this process to each of the seven input microphones and concatenate all resulting warped audio signals along the channel dimension.

Similar to Richard et al. [2021], we aim to estimate a geometric warpfeld  $\rho_{1:T}$  based on the speed of sound  $v_{\text{sound}}$  and the propagation distance difference between sound sources, input, and target microphones. We consider the aforementioned six body joints to be the potential sound sources  $\mathbf{p}_{1:S}^{(\text{src})}$  and simultaneously refer to the nose keypoint as the position of any input microphone  $\mathbf{p}_{1:S}^{(\text{in})}$ . Besides, the target microphone position  $(\theta, \phi)$  will be mapped into the Cartesian space  $\mathbf{x} = (x, y, z)$ . For a specific time stamp  $t \in \{1, \dots, T\}$ , we first identify the corresponding visual frame  $s = \lfloor t/1600 \rfloor$ . Accordingly, the time delay can be described as

$$\Delta t = (d_2 - d_1)/v_{\text{sound}} \quad \text{with} \quad d_1 = \|\mathbf{x} - \mathbf{p}_s^{(\text{src})}\|_2, d_2 = \|\mathbf{p}_s^{(\text{in})} - \mathbf{p}_s^{(\text{src})}\|_2. \quad (1)$$

In other words, we compute the distance  $d_1$  between the potential sound source and target microphone and the distance  $d_2$  between the potential sound source and input microphone. The time delay to

---

\*Work done during internship at Meta Reality Labs Research, Pittsburgh, PA, USA.

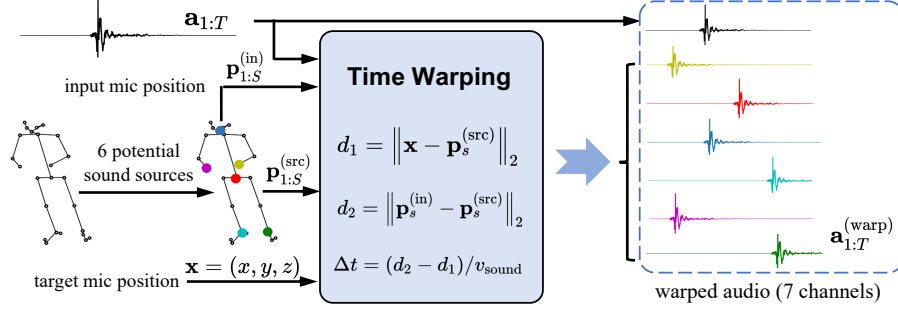


Figure 1: Overview of time warping. We manually select 6 keypoints as the locations of potential sound sources  $\mathbf{p}_{1:S}^{(\text{src})}$  and regard the nose keypoint as the input microphone position  $\mathbf{p}_{1:S}^{(\text{in})}$ . For a specific audio sample  $a_t$ , the corresponding visual frame  $s$  is first located, and the time warping is conducted according to the sound arriving time difference  $\Delta t$  between input and target microphones. Keeping the original input audio additionally to the warped signals, we obtain a set of audio  $\mathbf{a}_{1:T}^{(\text{warp})}$  containing seven channels (one original plus six-time warps) for each input microphone.

compensate for is determined by the distance between source and target minus the distance between source and input microphone.

Considering the audio is sampled at 48kHz, we represent the warpfeld as

$$\rho_t = \max(\rho_{t-1}, t - \Delta t \cdot 48000), \quad (2)$$

where monotonicity is ensured by  $\max(\rho_{t-1}, \cdot)$  following Richard et al. [2021]. It’s noteworthy that the causality property introduced in Richard et al. [2021] doesn’t hold here since the target microphone on the spherical array may receive the emitted sound before the head-mounted microphones. For example, a participant might stretch out their hand and make a snapping sound near the target microphone, but an arm-length away from the input microphone. With the predicted warpfeld  $\rho_t$ , the warped signal can be computed following Richard et al. [2021] as

$$a_t^{(\text{warp})} = (\lceil \rho_t \rceil - \rho_t) \cdot a_{\lfloor \rho_t \rfloor} + (\rho_t - \lfloor \rho_t \rfloor) \cdot a_{\lceil \rho_t \rceil}. \quad (3)$$

The time warping process is implemented on six selected sound sources, thus leading to six channels of warped audio. Additionally, we empirically find the original input audio  $\mathbf{a}_{1:T}$  can still facilitate the training and therefore keep it as an additional channel of the warped audio. Finally, the original input audio  $\mathbf{a}_{1:T}$  of a single input microphone is warped to 7-channel audio  $\mathbf{a}_{1:T}^{(\text{warp})}$  after time warping. We apply this time warping procedure to all input microphones and concatenate the results along the channel dimension.

## B Sound Field Encoding

In Section 4.4 we described how to decode the harmonic sound field coefficients into time domain signals at arbitrary locations in space. In this section we address the inverse problem: how to obtain the harmonic coefficients using signals captured by the  $N$  microphones arranged on a sphere around the center of the capture stage. We note that, due to spatial aliasing, the maximum harmonic order  $K$  of this encoding is limited by the number of microphones to  $K = \lfloor \sqrt{N} \rfloor - 1$ .

The signal  $s_i(t)$  captured by the  $i$ -th microphone located at  $(r_i, \theta_i, \phi_i)$  can be expressed in terms of sound field coefficients  $\beta_{nm}(\tau, f)$  as

$$\text{STFT}(s_i(t)) = \sum_{n=0}^K \sum_{m=-n}^n \beta_{nm}(\tau, f) h_n(kr_i) Y_{nm}(\theta_i, \phi_i), \quad (4)$$

where STFT denotes the short-time Fourier transform,  $\tau$  and  $f$  are, respectively, time and frequency bins,  $k = 2\pi f / v_{\text{sound}}$  is the wave number,  $v_{\text{sound}}$  is the speed of sound;  $Y_{nm}(\theta, \phi)$  represents the spherical harmonic of order  $n$  and degree  $m$  (angular spatial component), and  $h_n(kr)$  is  $n$ th-order spherical Hankel function (radial spatial component).

Given the STFT of all  $N$  microphone signals  $\mathbf{S}(\tau, f) = \text{STFT}([s_1(t), \dots, s_N(t)])^T$ , we can write Equation (4) as a linear system

$$\begin{aligned} \mathbf{S}(\tau, f) &= \begin{bmatrix} h_0(kr_1)Y_{00}(\theta_1, \phi_1) & \dots & h_K(kr_1)Y_{KK}(\theta_1, \phi_1) \\ \vdots & \ddots & \vdots \\ h_0(kr_N)Y_{00}(\theta_N, \phi_N) & \dots & h_K(kr_N)Y_{KK}(\theta_N, \phi_N) \end{bmatrix} \begin{bmatrix} \beta_{00}(\tau, f) \\ \vdots \\ \beta_{KK}(\tau, f) \end{bmatrix} \\ &= \mathbf{T}(f)\boldsymbol{\beta}(\tau, f), \end{aligned}$$

where  $\boldsymbol{\beta}(\tau, f)$  is a vector containing  $(K + 1)^2$  sound field coefficients, and  $\mathbf{T}(f)$  is a  $N \times (K + 1)^2$  matrix that links the coefficients to observations  $\mathbf{S}(\tau, f)$ . The sound field coefficients can then be obtained as

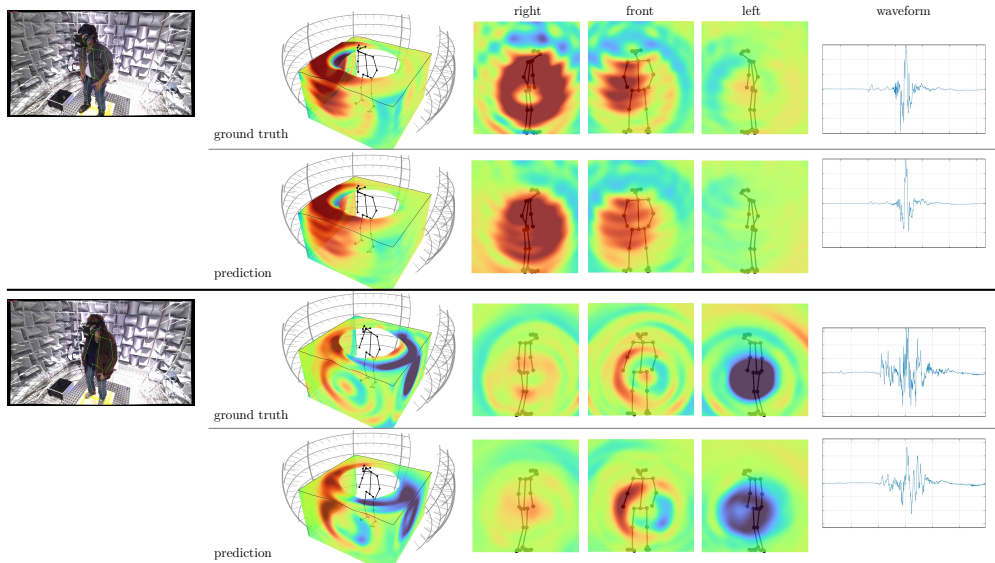
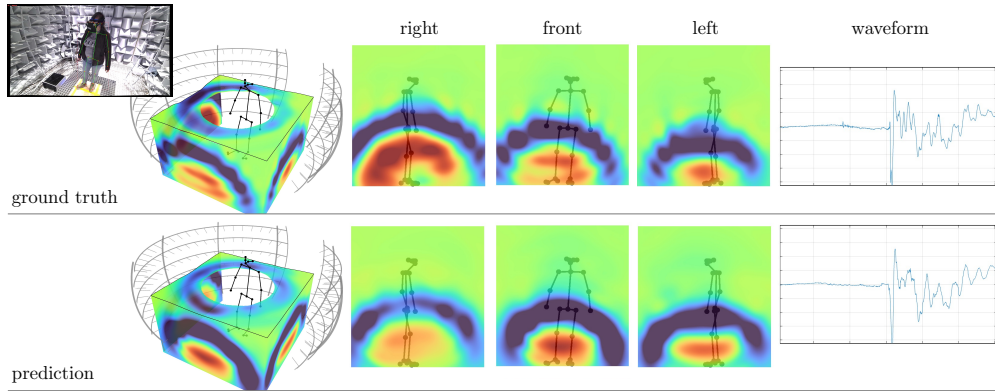
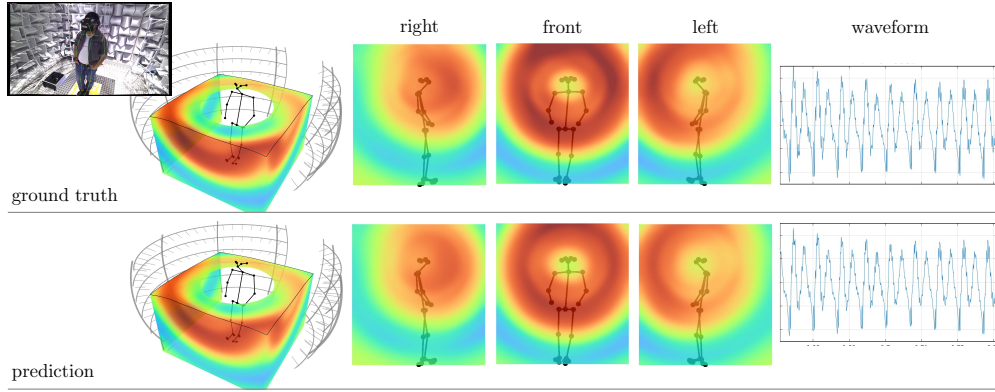
$$\boldsymbol{\beta}(\tau, f) = \mathbf{T}(f)^\dagger \mathbf{S}(\tau, f), \quad (5)$$

with  $(\cdot)^\dagger$  denoting the pseudo-inverse. For more details about spherical harmonic representation of sound fields and 3D recording of large areas see Williams [1999], Samarasinghe and Abhayapala [2012].

## C Audio-Visual Time Synchronization

The audio system is a distributed array of MADI-based analog-to-digital converters. Each analog-to-digital converter is phase-locked to each other via the same WordClock signal. The video system is a distributed array of Kinect Azure cameras. Each Kinect is set to “subordinate mode” where the unit will only take an image when they get a digital high signal from the sync-in port. The Kinect array is in a “star” pattern where one master unit sends a digital high signal out to a distribution amplifier. The master Kinect digital high signal is then propagated to all of the subordinate Kinects. We also record the digital high signal into our analog-to-digital audio converters so we have correspondence between a visual frame and an audio sample.

## D Additional examples





## E Failure case examples

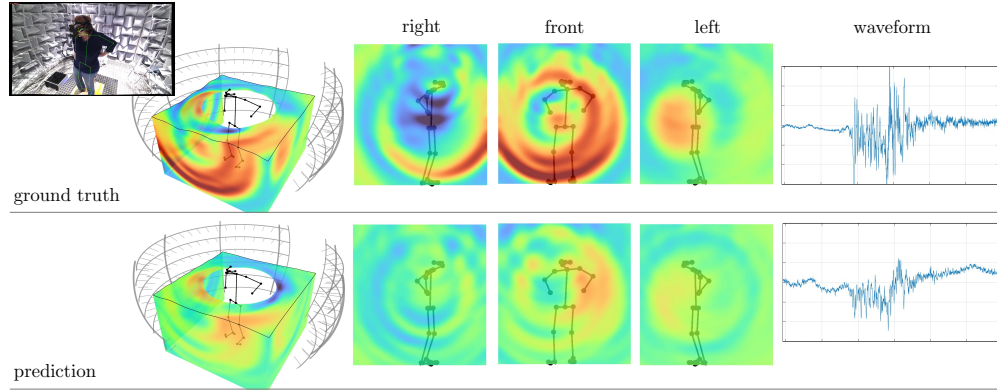


Figure 5: Not always predicted waveforms match the ground truth well. In this example, while sound source is reasonably localized, the signal energy is underestimated.

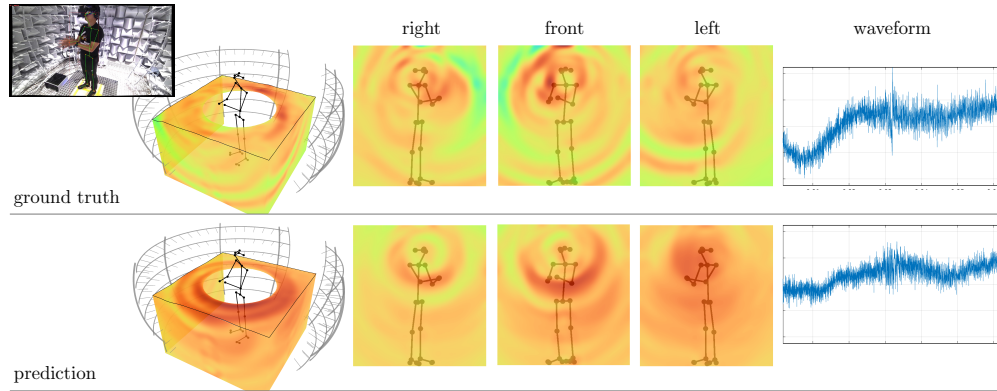


Figure 6: Hands rubbing example. Given noise-like nature and a very low energy of the signal, the model struggles to learn to correctly spatialize it.

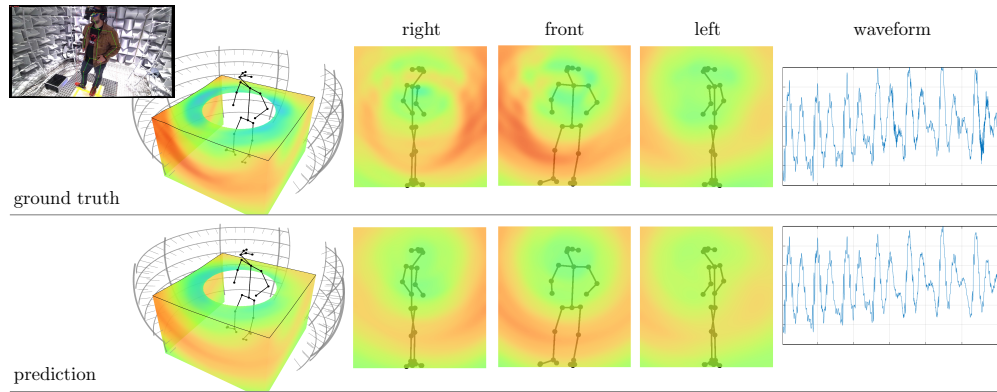


Figure 7: In this example a body tapping sound occurs at the same time as speech. Given disparity between signal energies, the model spatializes speech well but not the tapping sound.

Table 1: **Capture Script.** Each participant performed the following activities three times: once standing with light clothing (T-Shirt, light sweater, ...), once standing with heavy clothing that produces more noise (jackets, coats, ...), and once with light clothing while sitting.

activity	description	duration
<b>speech</b>		
<i>conversation_speech</i>	conversational speech	300s
<i>rainbow_L</i>	reading rainbow passage w/ left hand covering mouth	150s
<i>rainbow_R</i>	reading rainbow passage w/ right hand covering mouth	150s
<b>non-speech</b>		
<i>head_scratch</i>	scratch head w/ one hand at a time	15s
<i>face_rub</i>	play and rub face w/ one hand at a time	15s
<i>neck_scratch</i>	scratching neck w/ one hand at a time	15s
<i>back_scratch</i>	scratching back w/ one hand at a time	15s
<i>arms_stretching</i>	arms/fingers stretching	15s
<i>hands_rubbing</i>	rubbing hands ( <i>e.g.</i> , washing hands)	15s
<i>hands_arms_swipe</i>	swipe hands and arms	15s
<i>hands_scratch</i>	scratching arms and hands w/ one hand at a time	15s
<i>hands_punch</i>	punch other hand w/ one hand at a time	15s
<i>fist_bump</i>	fist bump	15s
<i>applause</i>	applause, <i>i.e.</i> clap hands in front of body	10s
<i>clapping</i>	clapping, covering as many spatial positions as possible	20s
<i>snapping_L</i>	snapping w/ left hand, covering as many spatial positions as possible	20s
<i>snapping_R</i>	snapping w/ right hand, covering as many spatial positions as possible	20s
<i>snapping_LR</i>	snapping w/ both hands, covering as many spatial positions as possible	15s
<i>trunk_twist</i>	twisting the trunk	10s
<i>body_stretching</i>	stretch body and yawn	10s
<i>chest_tap</i>	tapping chest w/ one hand at a time	20s
<i>belly_rub</i>	rubbing belly w/ one hand at a time	15s
<i>belly_tap</i>	tapping belly w/ one hand at a time	15s
<i>thighs_rub</i>	rubbing/scratching thighs w/ one hand at a time	15s
<i>thighs_tap</i>	tap thighs w/ one hand at a time	20s
<i>thighs_tap_LR</i>	tap thighs w/ both hands simultaneously	15s
<i>knees_rub</i>	rubbing/scratching behind knees w/ one hand at a time	15s
<i>body_itching</i>	itching and scratching the whole body	15s
<i>body_tapping</i>	tapping whole body (searching for phone, keys, wallet, ...)	15s
<i>feet_tap</i>	tap feet, one leg at a time	15s
<i>walk</i>	walk in place	15s
<b>mixed</b>		
<i>conversation_body</i>	conversational speech w/ body sounds and gestures	300s

## F Data Capture

Eight participants perform the script in Table 1 in three different settings: once standing with light clothing (T-Shirt, light sweater, ...), once standing with heavy clothing that produces more noise (jackets, coats, ...), and once with light clothing while sitting. During the capture, all participants wear a VR headset on which we mounted seven microphones. For the reading activities, participants are shown the text to read on the display of the VR headset. For non-speech activities, participants are shown an example video of an instructor performing the respective activity. For conversational tasks, participants are instructed to talk for the given duration about any topic of their choosing.

## References

- Alexander Richard, Dejan Markovic, Israel D Gebru, Steven Krenn, Gladstone Alexander Butler, Fernando Torre, and Yaser Sheikh. Neural synthesis of binaural speech from mono audio. In *International Conference on Learning Representations*, 2021.
- P. N. Samarasinghe and T. D. Abhayapala. 3D spatial soundfield recording over large regions. In *Proc. Int. Workshop on Acoustic Signal Enhancement (IWAENC)*, 2012.
- E. G. Williams. *Fourier Acoustics*. Academic Press, 1999.