

# Hearing the Room Through the Shape of the Drum: Modal-Guided Sound Recovery from Multi-Point Surface Vibrations: Supplementary materials

Shai Bagon    Matan Kichler    Mark Sheinin  
Weizmann Institute of Science, Israel

shai.bagon@weizmann.ac.il, matankic@gmail.com, mark.sheinin@weizmann.ac.il

## A. Spatially varying optical transfer and mode shape estimation

In Sec. 3.2 of the main manuscript, we simplified the relationship between the measured speckle shifts  $\mathbf{v}(\mathbf{x}_n, t)$  and the surface gradients  $\nabla_{\mathbf{x}} u(\mathbf{x}_n, t)$  by assuming a global optical transfer factor  $\beta$ . However, in a general imaging scenario, effects such as a high variance between the points' axial camera distances may yield different per-point scalars,  $\beta_n$ . In such circumstances, Eq. (7) becomes:

$$\mathbf{v}(\mathbf{x}_n, t) = \beta_n \nabla_{\mathbf{x}} u(\mathbf{x}_n, t). \quad (17)$$

Thus, substituting the modal expansion (Eq. (6)) into Eq. (17) now yields:

$$\begin{aligned} \mathbf{v}(\mathbf{x}_n, t) &= \gamma \beta_n \sum_{k=1}^K \nabla \phi_k(\mathbf{x}_n) (s(t) * g_k(t)) + \eta(\mathbf{x}_n, t) \\ &= \gamma \sum_{k=1}^K \nabla \phi_k^{\text{aug}}(\mathbf{x}_n) (s(t) * g_k(t)) + \eta(\mathbf{x}_n, t), \end{aligned} \quad (18)$$

where we define the *augmented* mode shape gradient as  $\nabla \phi_k^{\text{aug}}(\mathbf{x}_n) := \beta_n \nabla \phi_k(\mathbf{x}_n)$ .

Crucially, our procedure does not require explicitly decoupling the varying factors  $\beta_n$  from the true physical mode shape gradients  $\nabla \phi_k(\mathbf{x}_n)$ . Because our mode shape extraction is completely data-driven, it inherently absorbs this per-point scaling factor. Specifically, our numerical recovery of the mode shape gradients is extracted directly using Eq. (11) (reproduced below):

$$\nabla \hat{\phi}_k(\mathbf{x}_n) = \text{Re} \left\{ \frac{\mathbf{V}(\mathbf{x}_n, \hat{\omega}_k) \cdot \mathbf{V}_1(\mathbf{x}_0, \hat{\omega}_k)^*}{\mathbb{E}_{n,a} [|\mathbf{V}(\mathbf{x}_n, \hat{\omega}_k)|] \cdot |\mathbf{V}_1(\mathbf{x}_0, \hat{\omega}_k)|} \right\}. \quad (11)$$

Thus, in the case of non-negligible  $\beta_n$ , Eq. (11) directly approximates the augmented term  $\nabla \phi_k^{\text{aug}}(\mathbf{x}_n)$  up to a global normalization constant. The factors  $\beta_n$  will simply skew the numerically recovered mode shape gradients by a constant

per-point factor. Consequently, when we invert the model to solve for the recovered latent sound source  $\hat{s}^{\text{inv}}(t)$ , the optimization (Eq. (12)) remains completely agnostic to the true physical mode shape gradients of the surface:

$$\underset{s(t), \alpha_k}{\text{argmin}} \left\| \mathbf{v}(\mathbf{x}_n, t) - \sum_{k=1}^K \nabla \hat{\phi}_k^{\text{aug}}(\mathbf{x}_n) (s(t) * \hat{g}_k(t)) \right\|_2^2 + \lambda \|\hat{s}(t)\|_2^2. \quad (12)$$

Note that because the latent sound  $s(t)$  and the modal impulse responses  $\hat{g}_k(t)$  are convolved in our forward model, there is an inherent scale ambiguity. Any remaining global scaling discrepancies between the measured and physical domains are cleanly absorbed jointly by the optimized modal coupling coefficients  $\alpha_k$  (implicit within  $\hat{g}_k(t)$ ) and the overall amplitude of the recovered signal  $\hat{s}^{\text{inv}}(t)$ . Since  $s(t)$  is a dimensionless audio signal, this global scale factor simply alters the absolute playback volume without affecting the spectral content or acoustic fidelity. Thus, the precise physical calibration of  $\beta_n$  is bypassed without compromising the robust recovery of the acoustic signal.

## B. Full experimental results

We extensively evaluated our method on a variety of objects with diverse materials (wood, metal, plastic, rubber), geometries (planar, curved), and irregular shapes (such as the wooden binder and guitar body). The approach even works for solid objects like the Yoga block despite not being optimized for volumetric structures. Figure A shows all the objects along with the reconstructed spectrograms.

**Quantitative evaluation** We further evaluated the recovered audio using several established numerical scores. Note that we do not have a ground-truth target audio, only the reference source signal played by a speaker in front of each object. Hence, we quantitatively compare the recovered audio to the input source signal. We compare to three baselines: (i) *Single*: the raw vibrations measured at one point

ViSQOLAudio-NSIM $\uparrow$				
Experiment	Single	Avg	DnS	Ours
drum	0.20	0.19	0.18	<b>0.32</b>
frame	0.22	0.15	0.23	<b>0.46</b>
laptop	0.28	0.22	0.33	<b>0.44</b>
trash	0.19	0.24	0.24	<b>0.43</b>
guitar	0.16	0.26	0.24	<b>0.32</b>
binder	0.26	0.27	0.37	<b>0.43</b>
plate	0.19	0.17	0.23	<b>0.34</b>
drum (stereo)	0.26	0.26	0.27	<b>0.38</b>
yoga block	0.16	0.08	0.16	<b>0.29</b>
physio ball	0.35	0.35	0.39	<b>0.39</b>
balloon	0.27	0.17	0.39	<b>0.46</b>
<b>Mean:</b>	0.23	0.21	0.27	<b>0.39</b>

Table A. **ViSQOLAudio-NSIM** [2, 3] (higher is better). The raw perceptual similarity index used inside ViSQOLAudio before Mean Opinion Score (MOS) mapping. It measures structural similarity between cochleagram patches of the reference and degraded signals and is not tied to speech-specific MOS training. Higher values indicate closer perceptual similarity.

ViSQOLAudio-MOS $\uparrow$				
Experiment	Single	Avg	DnS	Ours
drum	<b>2.05</b>	1.90	1.96	1.63
frame	1.22	1.16	1.31	<b>1.96</b>
laptop	1.44	1.61	1.90	<b>2.06</b>
trash	1.07	1.12	1.05	<b>2.58</b>
guitar	<b>1.74</b>	1.59	1.66	1.58
binder	1.47	1.54	1.55	<b>2.45</b>
plate	1.24	1.43	1.45	<b>1.95</b>
drum (stereo)	1.69	1.65	1.55	<b>1.84</b>
yoga block	1.16	1.00	1.19	<b>2.67</b>
physio ball	1.37	1.54	1.32	<b>1.57</b>
balloon	1.58	1.43	<b>1.86</b>	1.81
<b>Mean:</b>	1.46	1.45	1.53	<b>2.01</b>

Table B. **ViSQOLAudio-MOS** [3] (higher is better). Perceptual audio-quality metric based on the ViSQOLAudio model. It compares a reference and degraded signal using a cochleagram front-end and patch-based similarity, then maps the result to a Mean Opinion Score (MOS) prediction. Originally trained on general-audio listening tests, it approximates subjective quality.

on the surface, (ii) *Avg*: naive averaging of all measured vibrations, and (iii) *DnS*: delay-and-sum where a single global temporal delay is estimated per point [1].

Table A reports ViSQOLAudio-NSIM [2, 3], which is the raw perceptual similarity index used inside ViSQOLAudio

Scale Invariant, Multi Resolution STFT distance $\downarrow$				
Experiment	Single	Avg	DnS	Ours
drum	4.68	4.21	4.60	<b>3.54</b>
frame	3.02	3.46	3.15	<b>1.87</b>
laptop	4.01	5.17	4.59	<b>3.88</b>
trash	3.65	3.27	3.46	<b>3.01</b>
guitar	3.23	<b>2.47</b>	2.59	3.49
binder	3.00	3.51	3.14	<b>2.72</b>
plate	3.22	3.60	3.82	<b>3.21</b>
drum (stereo)	<b>4.27</b>	4.56	4.41	4.32
yoga block	3.92	4.20	<b>3.66</b>	3.68
physio ball	<b>2.92</b>	3.40	3.16	4.00
balloon	<b>3.18</b>	3.39	3.96	3.73
<b>Mean:</b>	3.55	3.75	3.69	<b>3.40</b>

Table C. **Scale Invariant, Multi Resolution STFT distance** [4, 5] (lower is better). A distance metric between recovered signal and reference source signal. The measure is a scale invariant, multi resolution STFT distance.

Scale Invariant, Multi Resolution STFT distance (Perceptual weighting) $\downarrow$				
Experiment	Single	Avg	DnS	Ours
drum	3.65	3.39	<b>3.28</b>	3.39
frame	3.64	5.01	3.94	<b>2.01</b>
laptop	3.11	3.66	3.92	<b>2.91</b>
trash	3.56	3.57	3.88	<b>2.61</b>
guitar	3.45	2.91	3.04	<b>2.53</b>
binder	3.00	3.18	3.20	<b>2.48</b>
plate	3.44	3.57	4.37	<b>2.70</b>
drum (stereo)	3.67	<b>3.59</b>	4.16	3.97
yoga block	4.14	4.01	3.87	<b>3.68</b>
physio ball	<b>3.08</b>	3.22	3.12	3.35
balloon	<b>3.00</b>	3.12	3.42	3.58
<b>Mean:</b>	3.43	3.57	3.65	<b>3.02</b>

Table D. **Scale Invariant, Multi Resolution STFT distance (Perceptual weighting)** [4, 5] (lower is better). A distance metric between recovered signal and reference source signal. The measure is a scale invariant, multi resolution STFT distance with perceptual weighting.

before mean opinion score (MOS) mapping. Higher ViSQOLAudio-NSIM score indicates better similarity between the recovered audio and the input source signal. Our recoveries consistently score higher than all baselines. We further used the mean opinion score (MOS) mapping on top of ViSQOLAudio (Tab. B).

Inspired by [4, 5], We further used scale-invariant multi-

resolution distance computed in the space of short-time Fourier Transform (STFT). Table C shows the basic distance, while Tab. D shows a perceptually-weighted distance computed in this space. On average, our results have smaller distance to the reference input signal compared to all other baselines quantitatively showing the superiority of our audio recovery.

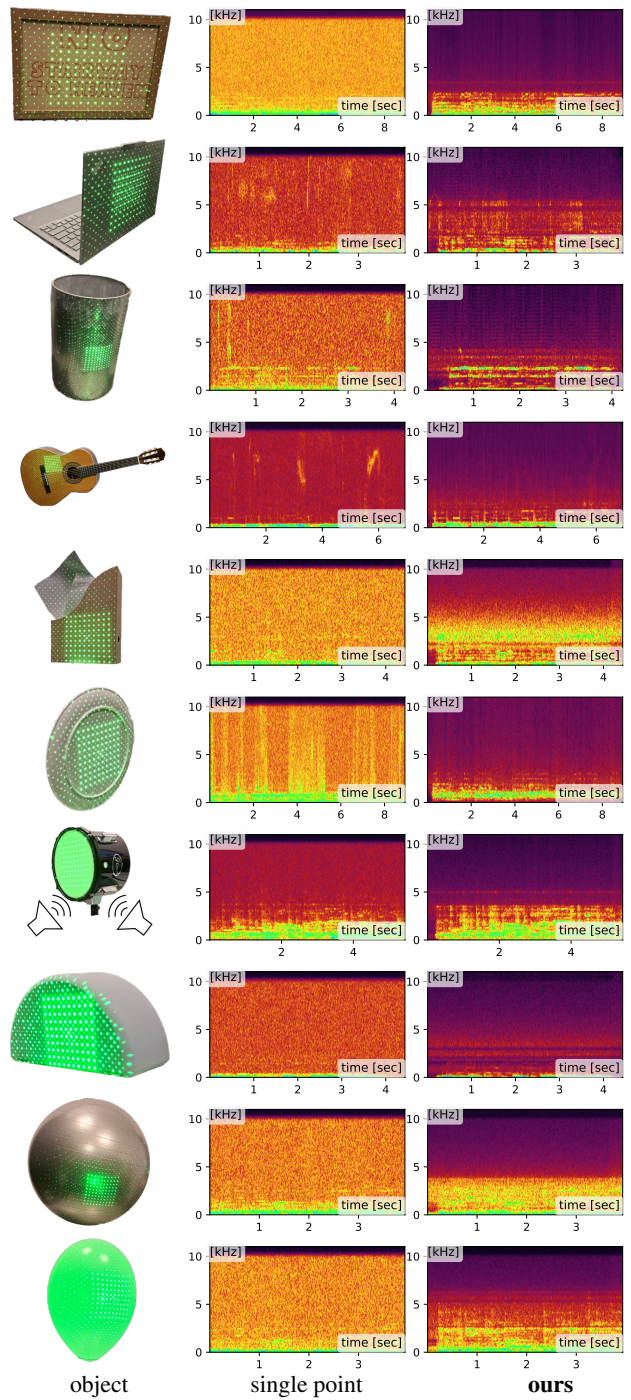


Figure A. Results across objects having various geometries and materials. All objects present challenges due to diverse materials, thicknesses, and shapes. Despite limited or irregular surface coverage, our method yields denoised reconstructions.

## References

- [1] Jacob Benesty, Israel Cohen, and Jingdong Chen. *Fundamentals of Signal Enhancement and Array Signal Processing*. John Wiley & Sons Singapore Pte. Ltd., 2017. 2
- [2] Andrew Hines and Naomi Harte. Speech intelligibility prediction using a neurogram similarity index measure. *Speech Communication*, 54(2):306–320, 2012. 2
- [3] Andrew Hines, Eoin Gillen, Damien Kelly, Jan Skoglund, Anil Kokaram, and Naomi Harte. ViSQOLAudio: An objective audio quality metric for low bitrate codecs. *The Journal of the Acoustical Society of America*, 137(6):EL449–EL455, 2015. 2
- [4] Christian J. Steinmetz and Joshua D. Reiss. auraloss: Audio focused loss functions in PyTorch. In *Digital Music Research Network One-day Workshop (DMRN+15)*, 2020. 2
- [5] Ryuichi Yamamoto, Eunwoo Song, and Jae-Min Kim. Parallel WaveGAN: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram. In *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2020. 2