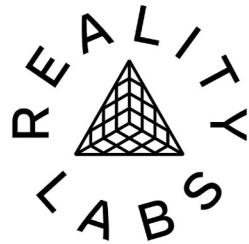




University
of Rochester

Meta



UNIVERSITY OF
MARYLAND

Towards Perception-Informed Latent HRTF Representations

You (Neil) Zhang^{1,2}, **Andrew Francl**², **Ruohan Gao**³, **Paul Calamia**²,
Zhiyao Duan¹, **Ishwarya Ananthabhotla**²

IEEE
WASPAA
2025

¹ University of Rochester

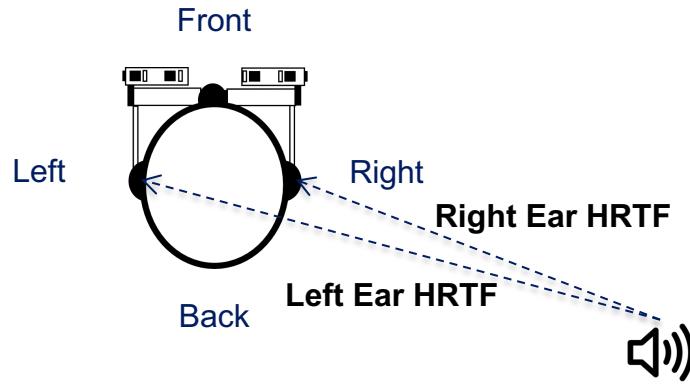
² Meta Reality Labs Research

³ University of Maryland

Tahoe City, CA - Oct 13, 2025

Head-Related Transfer Function (HRTF)

HRTF models the **acoustic filtering** effect of a listener's **head, ears, and torso** to enable 3D sound localization.



Left ear HRTF magnitudes (dB) of the midsagittal plane of one subject

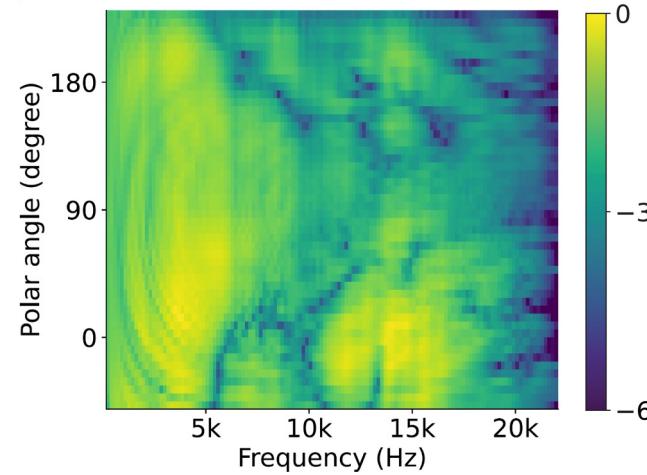


Figure from [Zhang+2023]

HRTF is unique to each person due to differences in ear, head, and torso shape.

HRTF Applications: Virtual Spatial Audio Rendering

HRTFs encode human spatial cues to deliver immersive 3D sound.



Headphones



AR smart glasses



VR headsets

Measure HRTFs

- An anechoic room
- Multiple loudspeakers on motorized arc
- Two microphones
- Head motion control

Time-consuming & Resource-intensive!

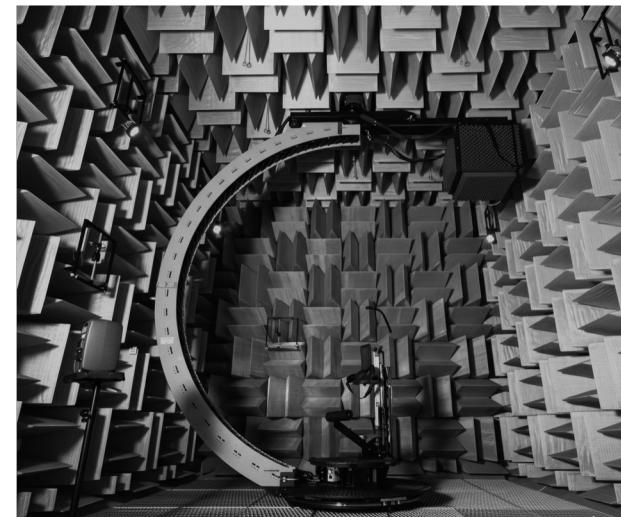
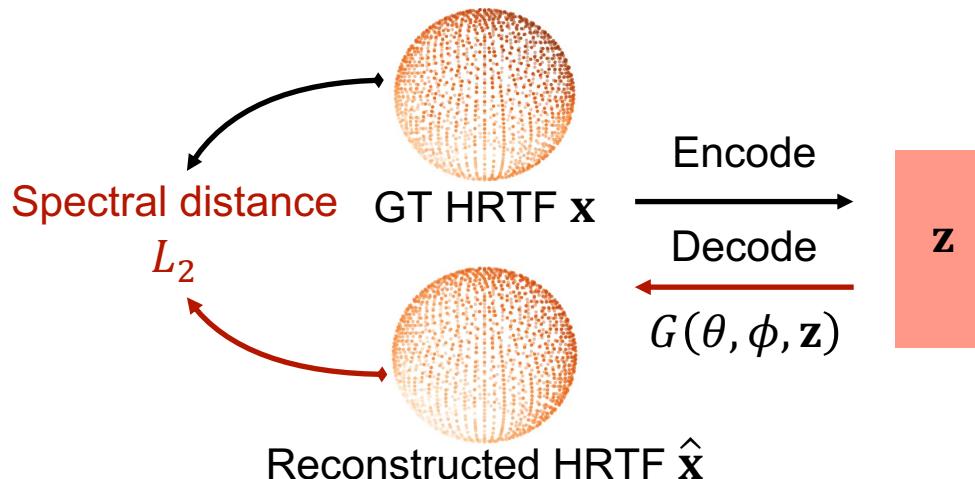


Figure from https://facebookresearch.github.io/SS2_HRTF/

Deep Learning for HRTF Representations

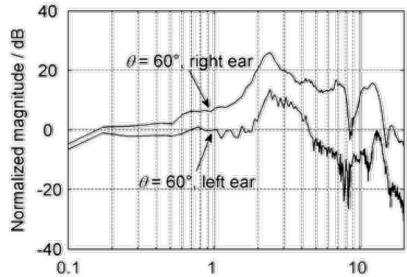
HRTF representation learning models

- Convolutional Autoencoder (CAE)
- Implicit Neural Representation (INR)

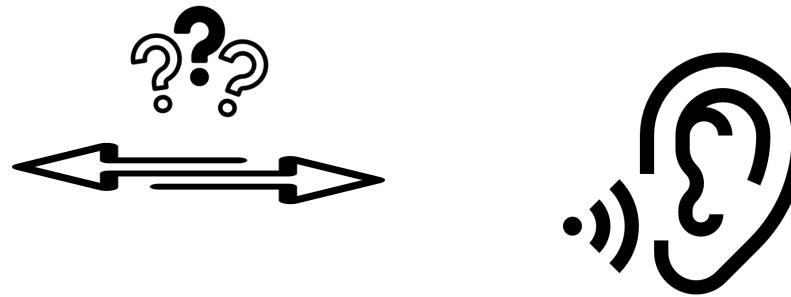


Motivation

Most **existing** models are trained and evaluated with **spectral reconstruction**.



Spectral reconstruction



Perceptually plausible HRTF

Our contributions

Goal: Learn HRTF representations that **more accurately reflect perceptual correlation**, to enable better HRTF personalization for unseen users



- We study how well **existing** latent HRTF representations preserve perceptual relations, and **introduce the benchmark** for evaluating this.
- We propose **a method for improving** on this benchmark.
- We demonstrate **practical utility** for HRTF personalization.

1. How well do **existing** learned HRTF representations **preserve perceptual relations?**

HRTF Perception

Perceptual benefits of your *personal HRTF*:

- Reduced **Coloration** (less unwanted spectral distortion)
- Improved **Externalization** (sound appears outside the head)
- Enhanced **Localization** (accurately placing sounds in 3D space)

Coloration

How do we mathematically model these?

Externalization

Localization

Computational Auditory Modeling

Coloration: Predicted Binaural Coloration [McKenzie+2022]

Externalization: Auditory Externalization Perception [Baumgartner&Majdak2021]

Localization: Difference of Root Mean Square Error in Polar Angles [Barumerli+2023]

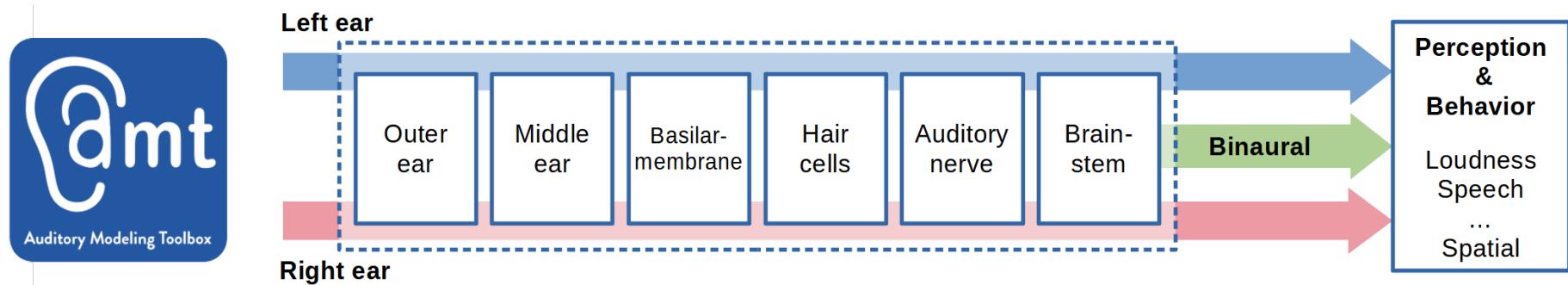


Figure from amttoolbox.org

McKenzie, Thomas, et al. "Predicting the colouration between binaural signals." *Applied Sciences* 2022.

Baumgartner, Robert, and Piotr Majdak. "Decision making in auditory externalization perception: model predictions for static conditions." *Acta Acustica* 2021 111

Barumerli, Roberto, et al. "A Bayesian model for human directional localization of broadband static sound sources." *Acta Acustica* 2023.

Computational Auditory Modeling

Coloration: **PBC** [McKenzie+2022]

Externalization: **AEP** [Baumgartner&Majdak2021]

Localization: **DRMSP** [Barumerli+2023]

Objective Perceptual Metrics

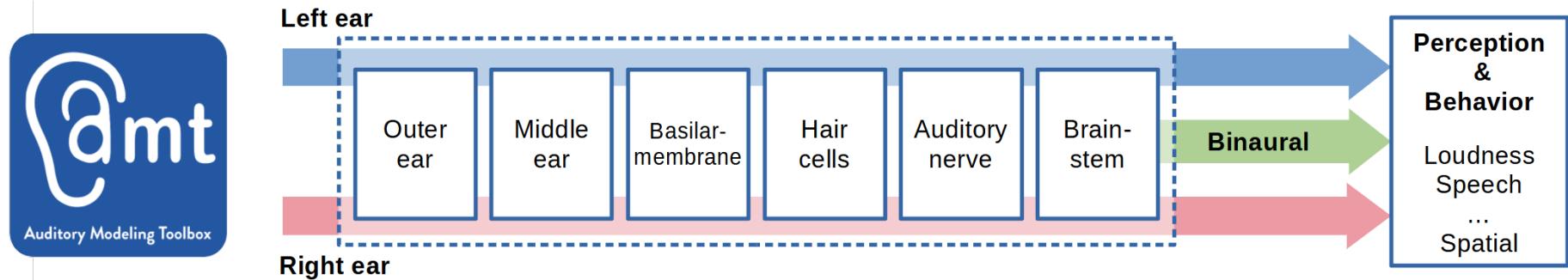


Figure from amttoolbox.org

McKenzie, Thomas, et al. "Predicting the colouration between binaural signals." *Applied Sciences* 2022.

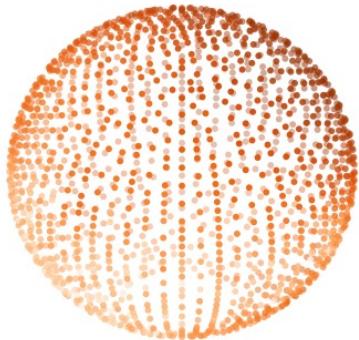
Baumgartner, Robert, and Piotr Majdak. "Decision making in auditory externalization perception: model predictions for static conditions." *Acta Acustica* 2021 12

Barumerli, Roberto, et al. "A Bayesian model for human directional localization of broadband static sound sources." *Acta Acustica* 2023.

Experimental Setup

SS2 HRTF Database

- 1625 measurement locations
- 48 kHz sampling rate
- 78 subjects (65 for training,
13 for testing)



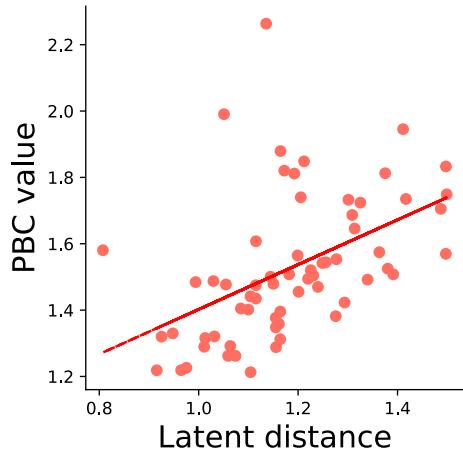
Train AI models with spectral reconstruction for HRTF data

Compute pairwise **latent** distance across subjects

Compute pairwise **perceptual** distance across subjects

Alignment Between Latent Space and Perceptual Metrics

Pearson correlation (pairwise latent distances vs. perceptual distances)



$$\rho_{A,B} = \frac{\mathbb{E}[(A - \mu_A)(B - \mu_B)]}{\sigma_A \sigma_B}$$

A higher positive correlation indicates better alignment with human perception.

Partitions	PBC	AEP	DRMSP
train	0.60	0.60	0.40

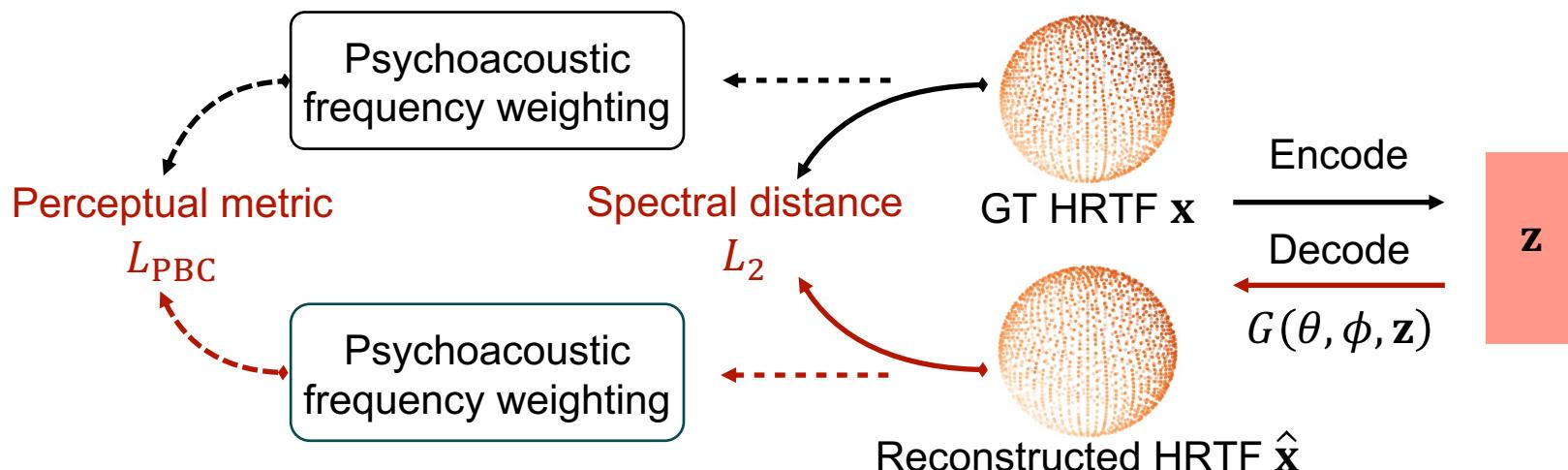
Minimizing spectral distances leads to limited perceptual correlation.

2. How do we **align** latent HRTF representations with **perception-informed space**?

Aligning with Perception-Informed Space

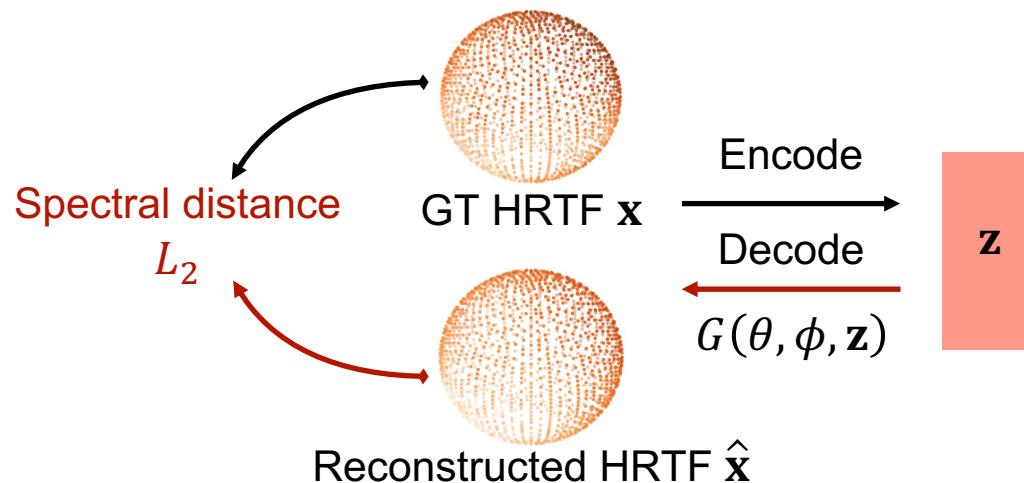
If the perceptual metric is differentiable, just add a straightforward perceptual loss.

- This only applies to PBC, which we reimplemented with PyTorch.



Aligning with Perception-Informed Space (Cont'd)

If the perceptual metric is not differentiable (AEP, DRMSP)



Metric multi-dimensional scaling (MMDs)

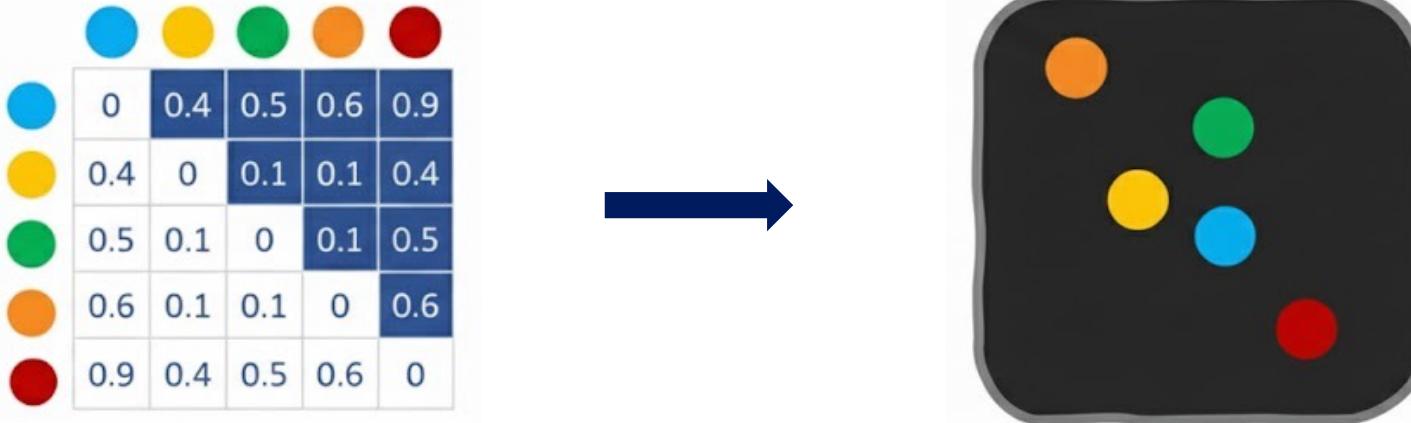
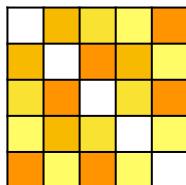


Figure from https://youtu.be/VKSJayDi_IQ post-processed by Gemini

Aligning with Perception-Informed Space (Cont'd)

If the perceptual metric is not differentiable (AEP, DRMSP)



Metric multi-dimensional scaling (MMDS)

Pairwise perceptual distance matrix \mathbf{M}

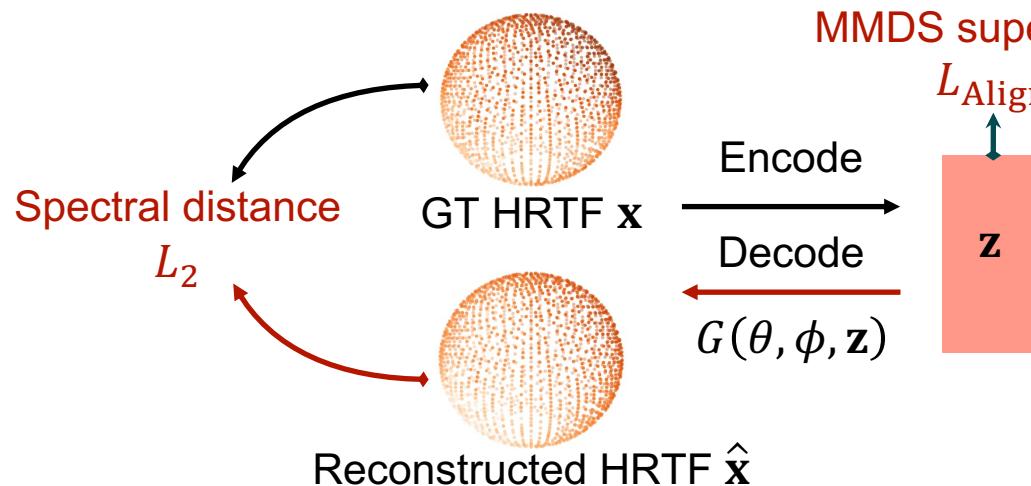


\mathbf{z}_{MDS}

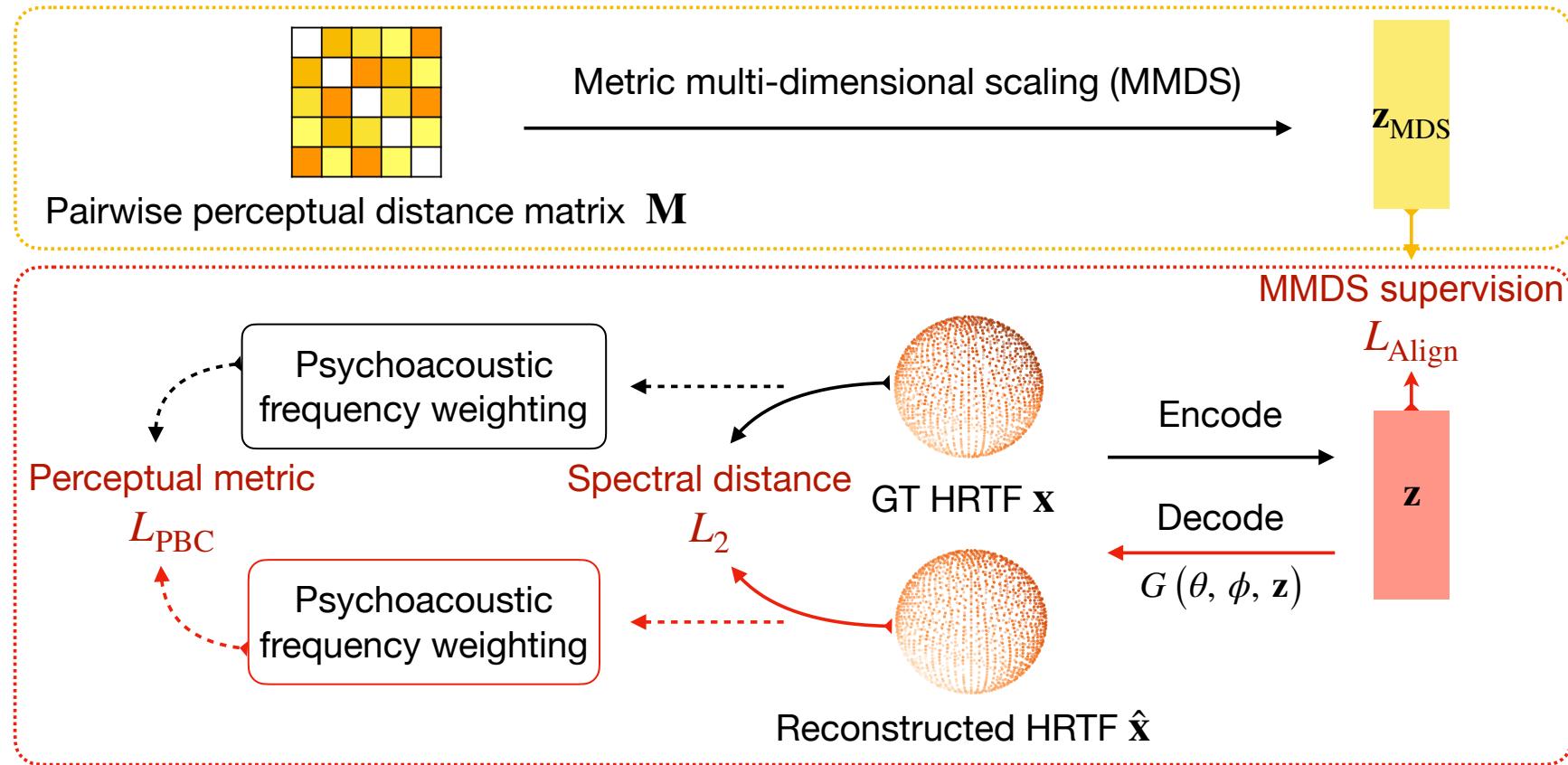
MMDS supervision

L_{Align}

This can also be applied to differentiable metrics.

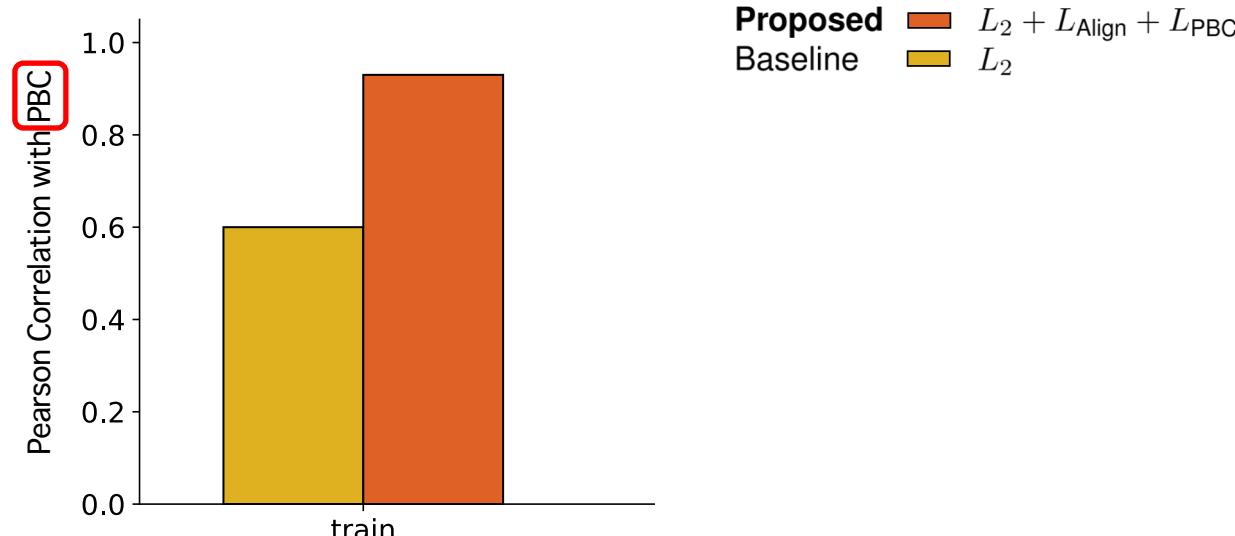


Overall Pipeline to Align Latent HRTF Representations



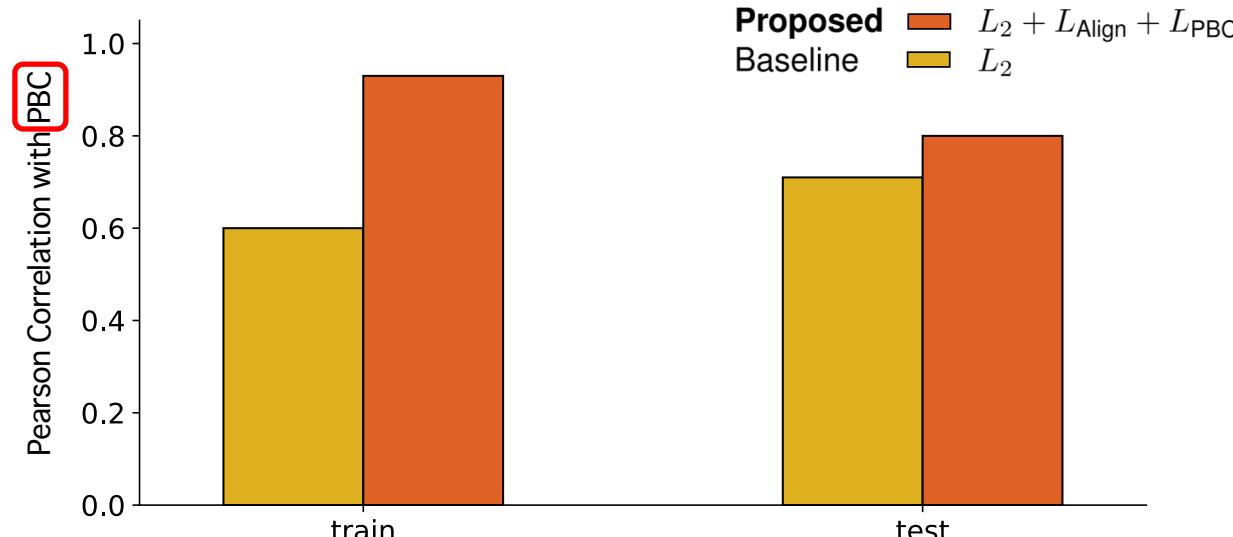
Results: Objective Perceptual Correlation Evaluation on PBC

- Our proposed method **achieves better alignment** with perception-informed space.
- The perceptual correlation learned in training transfer to test subjects (unseen).



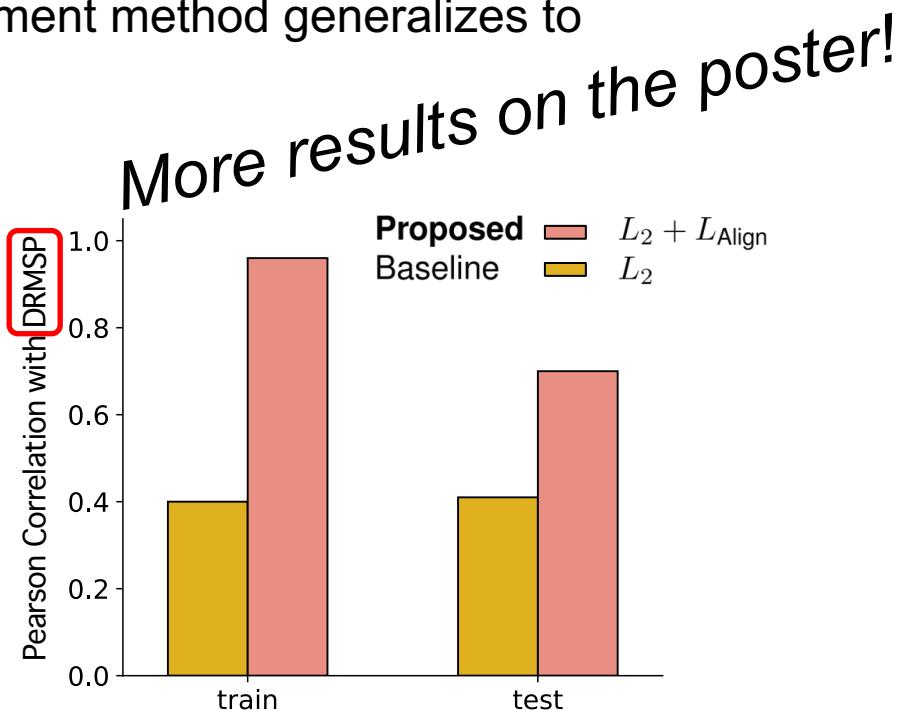
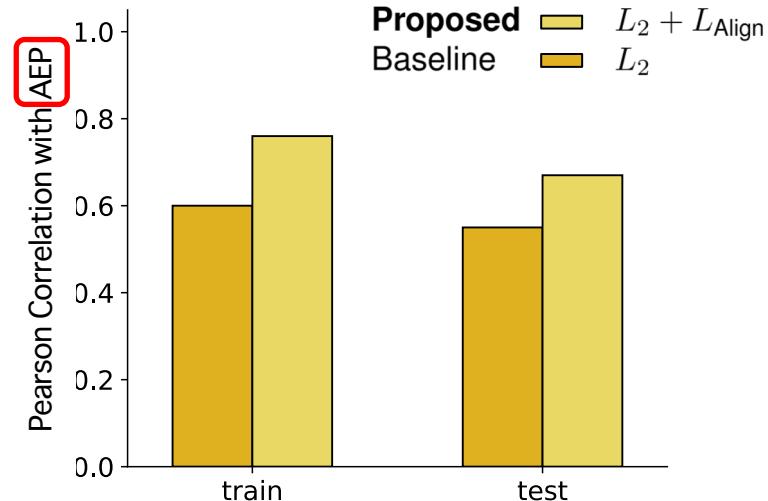
Results: Objective Perceptual Correlation Evaluation on PBC

- Our proposed method **achieves better alignment** with perception-informed space.
- The perceptual correlation learned in training **transfer to test subjects (unseen)**.



Generalization to AEP and DRMSP Metrics

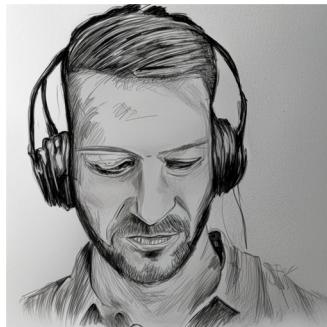
- YES, our proposed correlation improvement method generalizes to externalization and localization.



More results on the poster!

Application: Personalized HRTF Selection

For each of the 13 test (unseen) subjects, we select the **nearest** HRTFs from the 65 training subjects, based on the learned latent representations.



Methods	Best candidate	
	DRMSP↓	SDE (dB)↓
$L_2 + L_{\text{Align}}$	3.20	2.12
L_2	4.21	2.07

HRTFs selected by our proposed method yield
lower perceptual distances with slightly higher SDE.

More results on the poster!

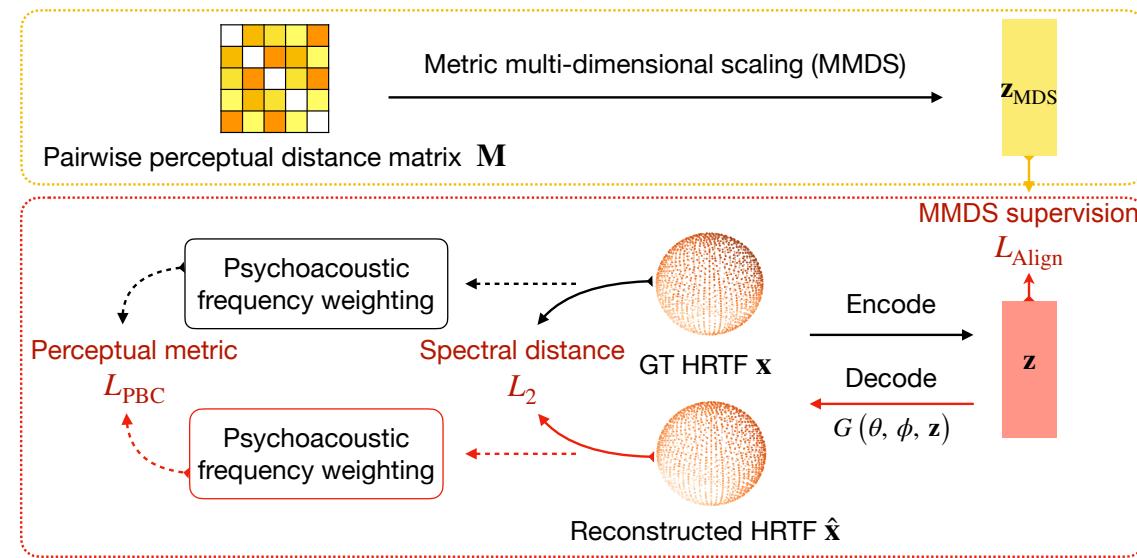
Limitations and Future Work

Limitations

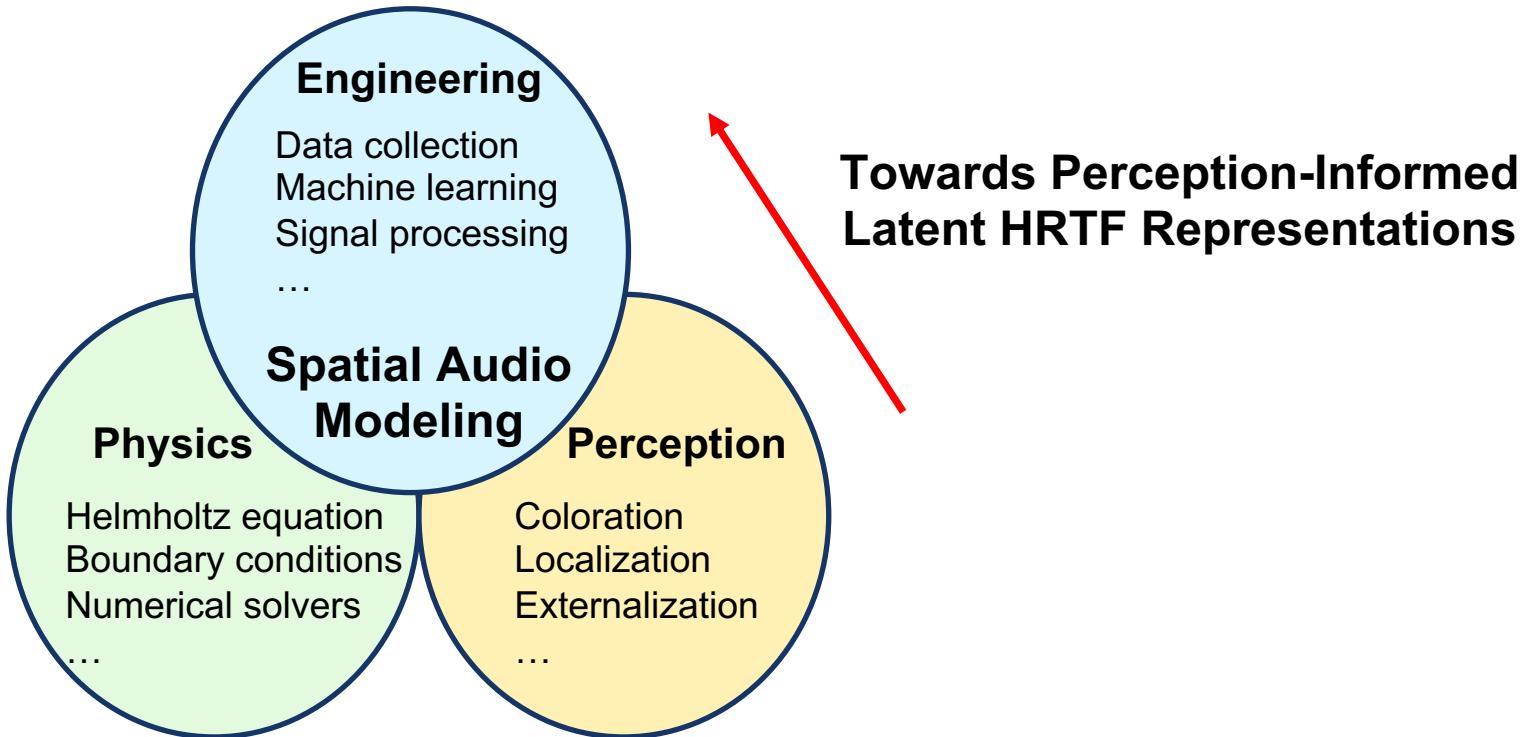
- Objective metrics vs. **listening experience**
- MMDS assumes symmetric dissimilarity
- Ignoring phase information

Future Work

- Subjective validation
- Extension to binaural synthesis



Future Vision



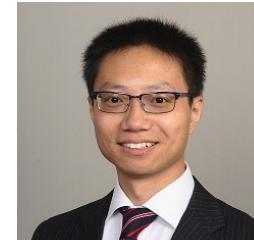
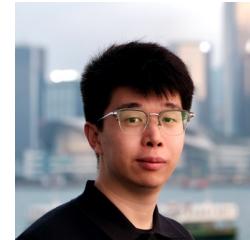
Takeaways

HRTF representation isn't just about spectral error —
perceptual distance is equally important!

Our proposed recipe: MMDS supervision + perceptual loss

Thank you! Questions?

Check out our paper and poster
for more, and let's discuss more!



<Backup Slides>

Spectral distance

- Spectral Difference Error (SDE)

$$\text{SDE}_k(H, \hat{H}) = \frac{1}{L} \sum_{\theta, \phi} \left| 20 \cdot \log_{10} \left(\frac{H(\theta, \phi, k)}{\hat{H}(\theta, \phi, k)} \right) \right|$$

The diagram illustrates the components of the SDE formula. It consists of four blue-bordered boxes connected by arrows to the corresponding parts of the equation:

- A box labeled "ground-truth linear-scale magnitude" has an arrow pointing to $H(\theta, \phi, k)$.
- A box labeled "# spatial locations" has an arrow pointing to L .
- A box labeled "predicted linear-scale magnitude" has an arrow pointing to $\hat{H}(\theta, \phi, k)$.
- A box labeled "frequency index" has an arrow pointing to k .

The median SDE across all frequency bins was computed to obtain a single SDE value.

SS2 HRTF Dataset

- High-resolution HRTF database with 1625 measurement locations
- 48 kHz sampling rate
- 78 subjects

