



Australian
National
University

Room Impulse Response Reconstruction Based on Spatio-temporal-spectral Features Learned from a Spherical Microphone Array Measurement

Amy Bastine¹, Thushara D. Abhayapala¹, Jihui (Aimee) Zhang^{2,1}

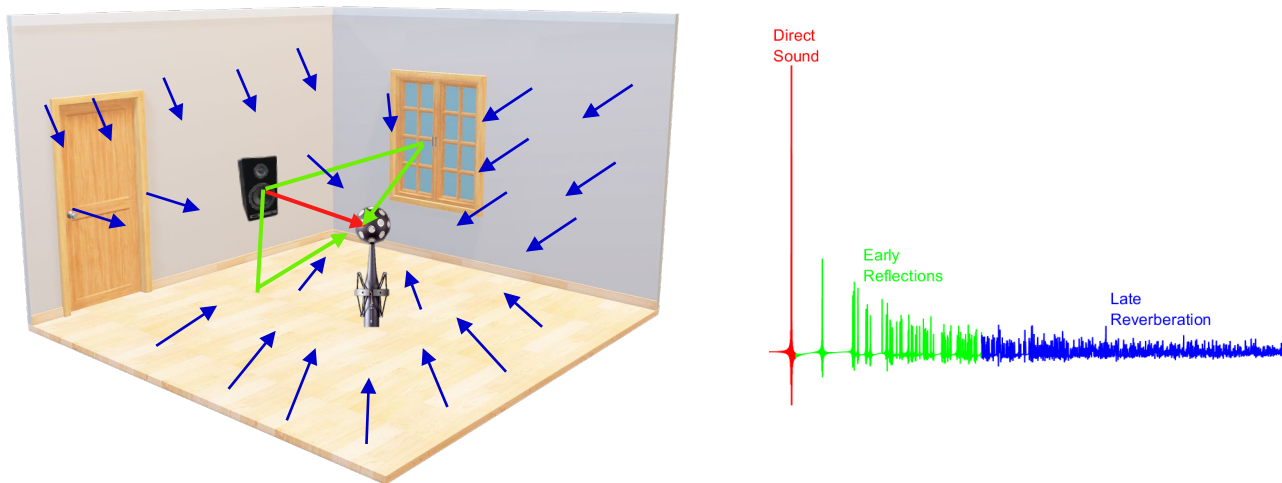
*¹Audio & Acoustic Signal Processing Group
College of Engineering, Computing and Cybernetics
The Australian National University
Canberra, Australia*

*²Signal Processing, Audio and Hearing Group
Institute of Sound and Vibration Research (ISVR)
University of Southampton
Southampton, United Kingdom*

2023 International Conference on Acoustics, Speech, and Signal Processing (ICASSP)
Special Session: Data Driven and Machine Learning based Room Acoustic Modeling

Room Impulse Responses (RIRs)

Fundamental representation of a room acoustic system



- Large-scale RIR measurements required to determine room's response to different source-listener configurations
- **RIR reconstruction methods**
 - Enable estimation of listener experience outside the measurement positions
 - Reduce measurement costs

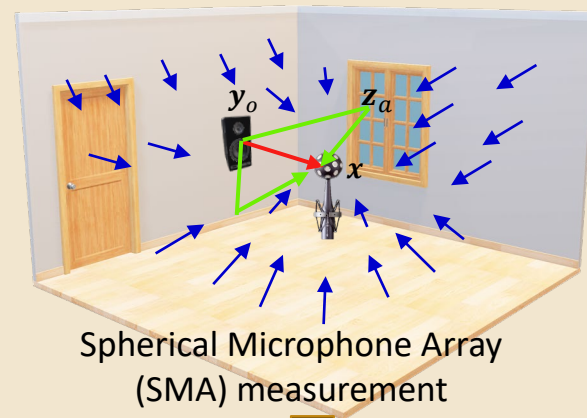
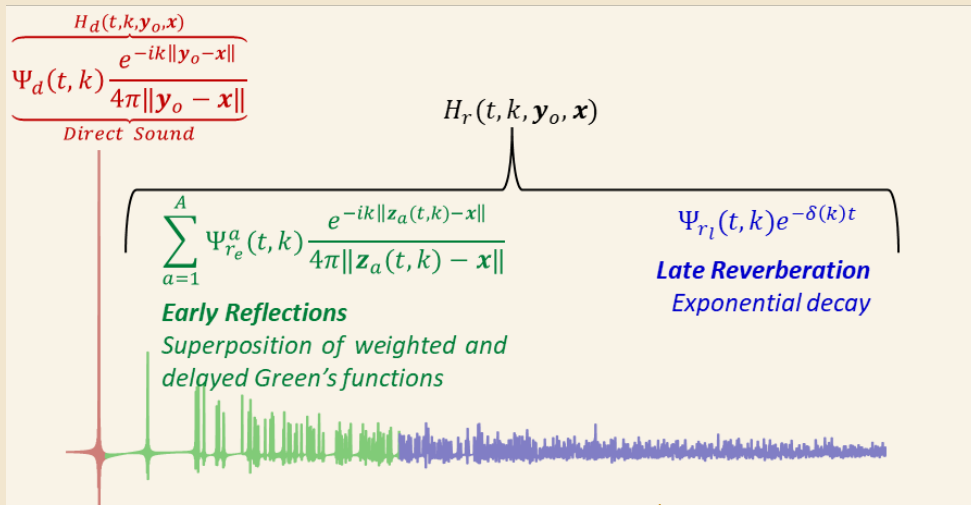
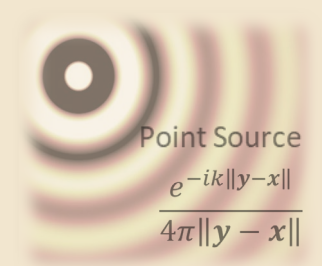
RIR Reconstruction

Existing Methods

- **Model-based methods**: Useful for prediction, but may not accurately emulate the intended room response
 - *Wave-based*: Computationally expensive
 - *Geometrical* models like Image-Source Method (ISM), Ray Tracing: Limited to high frequencies
 - *Hybrid Approaches*
- **Data-driven methods** based on existing measurements generate more authentic RIRs
 - *Machine learning*: Requires large amounts of training data
 - *Interpolation*: Requires distributed grid of microphone measurements
 - *Extrapolation-based Parametric* methods: Minimal measurement and computational cost
 - Spherical Microphone Arrays (SMAs) & Spherical Harmonics-based processing ⇒
Higher-order soundfield information ⇒ Improve reconstruction performance

RIR Reconstruction

Using spatio-temporal-spectral features learned from a SMA measurement



Spherical Microphone Array (SMA) measurement

Features of room reflections:

Reflection source locations: $z_a(t, k)$

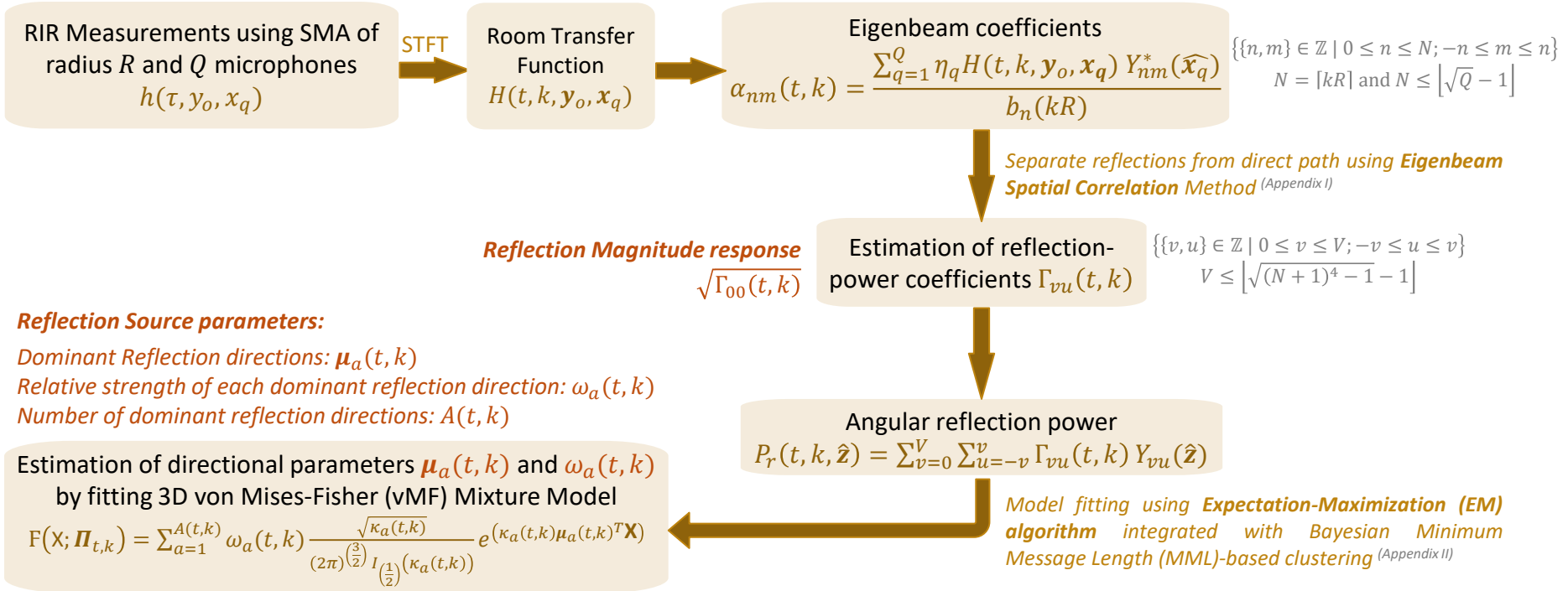
Reflection magnitude response: $\Psi_r(t, k)$

Parameter Estimation
Eigenbeam vMF-based
Room Acoustic Analyzer



RIR Reconstruction

1. Parameter estimation using *Eigenbeam vMF-based Room Acoustic Analyzer*

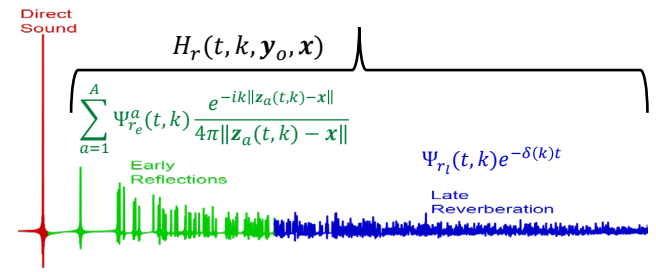


A. Bastine, T. D. Abhayapala, and J. Zhang, "Time-frequency-dependent directional analysis of room reflections using Eigenbeam processing and von Mises-Fisher clustering," *J. Acoust. Soc. Amer.*, vol. 151, no. 5, pp. 2916–2930, May 2022.



RIR Reconstruction

2. Synthesis of Reflection Transfer Function



$$\tilde{H}_r(t, k, \mathbf{y}_o, \mathbf{x}) = \begin{cases} \text{Early Reflections (superposition of weighted and delayed Green's functions)} & \text{for } \{t, k\} \text{ frames where } \mathbf{\Pi}_{t, k} \text{ exists} \\ \text{Late Reverberations (Exponential Decay)} & \text{; otherwise} \end{cases}$$

$$= \begin{cases} \sqrt{\Gamma_{00}(t, k)} \sum_{a=1}^{A(t, k)} A(t, k) \omega_a(t, k) \frac{e^{-ik\|z_a(t, k) - \mathbf{x}\|}}{4\pi\|z_a(t, k) - \mathbf{x}\|} e^{-ikd_t} & \text{; for } \{t, k\} \text{ frames where } \mathbf{\Pi}_{t, k} \text{ exists} \\ \sqrt{\Gamma_{00}(t, k)} \frac{\tilde{H}_r(t-1, k, \mathbf{y}_o, \mathbf{x})}{\sqrt{\Gamma_{00}(t-1, k)}} e^{-\delta(k)t_f} e^{-ikd_t} & \text{; otherwise} \end{cases}$$

$\sqrt{\Gamma_{00}(t, k)}$: Time-frequency-dependent magnitude response of room reflections

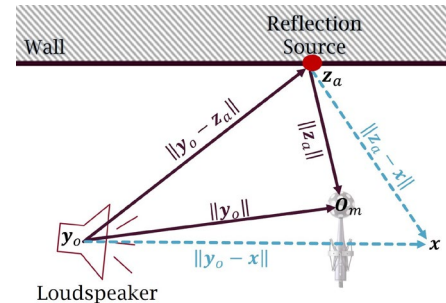
$A(t, k)\omega_a(t, k)$: Directional amplitude scaling

$\mathbf{z}_a(t, k)$: Location of room surface point in the direction of $\boldsymbol{\mu}_a(t, k)$ (Approximate room dimensions known)

e^{-ikd_t} : Phase-shift to align the response to the corresponding STFT time frame

d_t : $(t-1)ct_f$ with t_f being the time gap between STFT frames

$\delta(k)$: Decay rate calculated from $\sqrt{\Gamma_{00}(t, k)}$



RIR Reconstruction

3. Combining Direct and Reflection Components

$$\tilde{H}_d(t, k, \mathbf{y}_o, \mathbf{x}) = \Psi_d(t, k) \frac{e^{-ik\|\mathbf{y}_o - \mathbf{x}\|}}{4\pi\|\mathbf{y}_o - \mathbf{x}\|}$$

↓ ISTFT

$$\tilde{h}_d(\tau, \mathbf{y}_o, \mathbf{x})$$

↓ Normalize

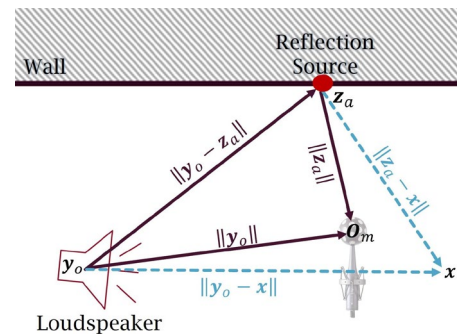
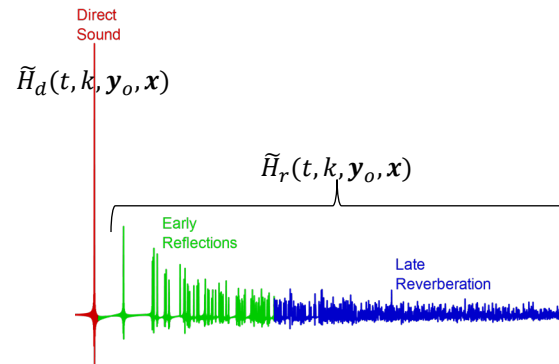
$$\tilde{h}(\tau, \mathbf{y}_o, \mathbf{x}) = \tilde{h}_d(\tau, \mathbf{y}_o, \mathbf{x}) + \tilde{h}_r\left(\tau - \frac{\|\mathbf{y}_o - \mathbf{z}_a(t=1, k)\|}{c}, \mathbf{y}_o, \mathbf{x}\right) e^{-\frac{\bar{\delta}d_p}{c}}$$

$$\tilde{H}_r(t, k, \mathbf{y}_o, \mathbf{x})$$

↓ ISTFT

$$\tilde{h}_r(\tau, \mathbf{y}_o, \mathbf{x})$$

↓ Normalize



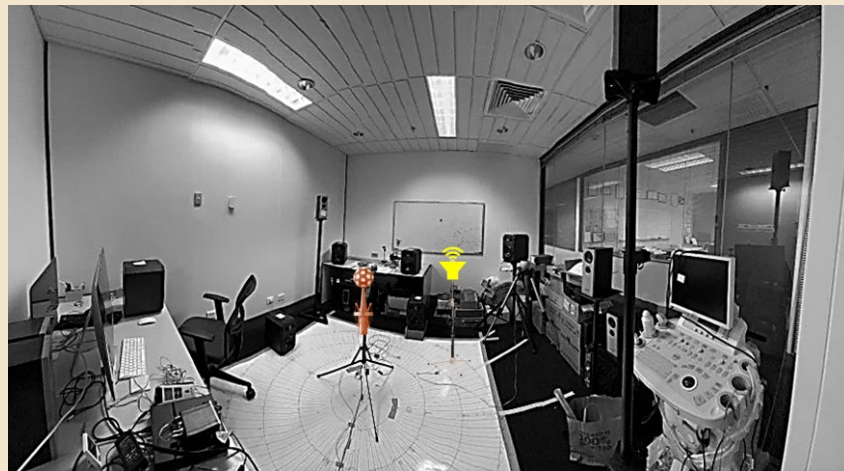
$$d_p = \|\mathbf{y}_o - \mathbf{z}_a(t=1, k)\| + \|\mathbf{z}_a(t=1, k) - \mathbf{x}\|$$

$$\bar{\delta} = \frac{3 \ln(10)}{T_{60}}$$

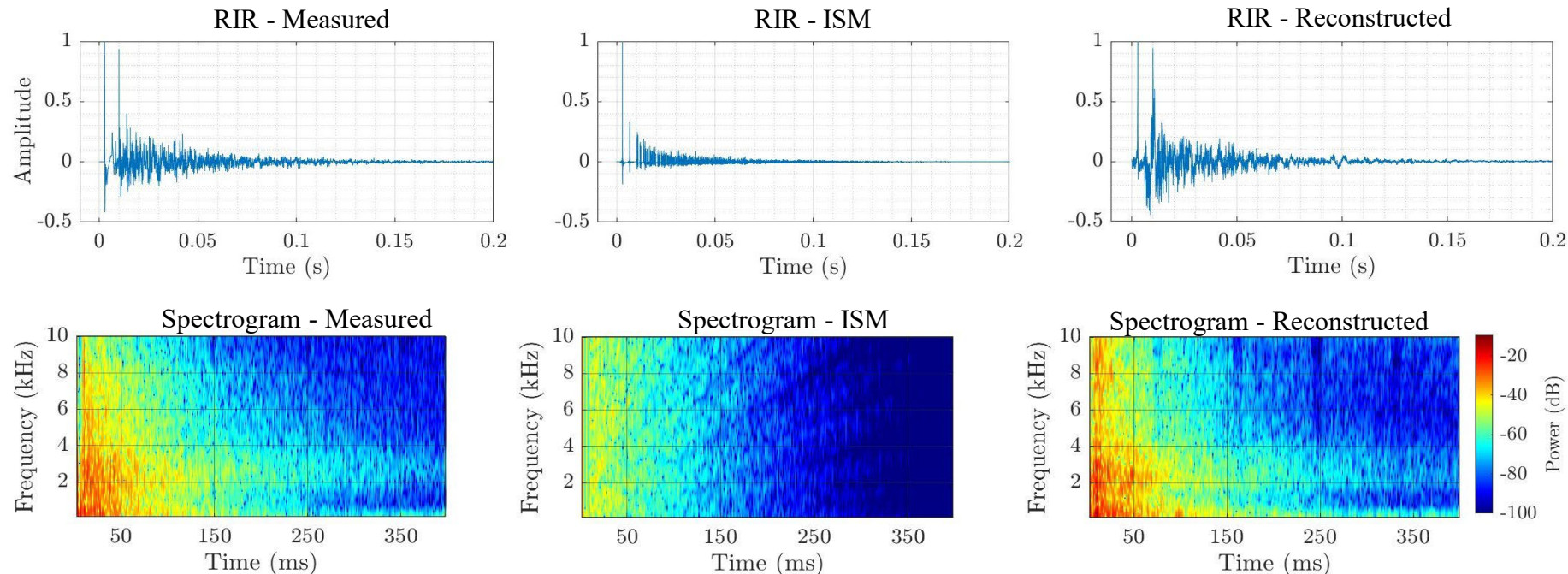


Experimental Analysis

- Room of size $3.54 \times 4.06 \times 2.7$ m and $T_{60} = 0.329$ s
- Source located at the spherical coordinate $\mathbf{y}_o = (1, 90^\circ, 40^\circ)$
- Parameters estimated from RIRs recorded by an EM32 Eigenmike (32-element rigid SMA with radius $R = 0.042$ m)
 - Maximum number of dominant reflection directions ($A(t, k)$) set to 5
 - Total of 7448 reflection sources $\mathbf{z}_\alpha(t, k)$ identified from all STFT frames
- Performance compared with measured RIRs and conventional Image Source Method (ISM)
 - ISM-based RIRs generated with maximum image order (Number of image sources in the order of 10^5)



Experimental Analysis

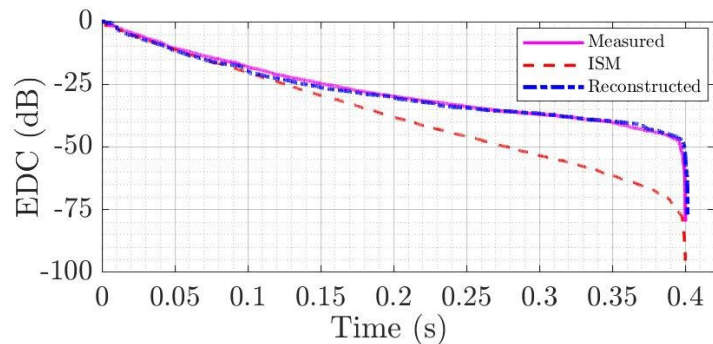


RIRs and corresponding spectrograms for a source located at $\mathbf{y}_o = (1, 90^\circ, 40^\circ)$ and receiver at $\mathbf{x} = (0.042, 69^\circ, 0^\circ)$

$$\text{NMSE (ISM, Measured)}_{0.4s} = -5 \text{ dB}$$

$$\text{NMSE (Reconstructed, Measured)}_{0.4s} = -17 \text{ dB}$$

Experimental Analysis



Energy Decay Curves (EDC) of RIRs obtained using normalized Schroeder integration method for a source located at $\mathbf{y}_o = (1, 90^\circ, 40^\circ)$ and receiver at $\mathbf{x} = (0.042, 69^\circ, 0^\circ)$.

Features of early and late reflections preserved in the reconstructed RIRs

Mean and Standard Deviation (STD) of objective room acoustic parameters calculated for the 32 receiver positions of EM32. $e_{I,M}$ and $e_{R,M}$ represent deviation errors of the mean parameter values of ISM-based and Reconstructed RIRs from the measured RIR, respectively.

	Parameters	Measured			ISM			Reconstructed		
		Mean	STD	JND	Mean	STD	$e_{I,M}$	Mean	STD	$e_{R,M}$
Early Decay Time	EDT (s)	0.24	0.007	5% = 0.012	0.14	0.001	0.1	0.23	0.0016	0.01
Clarity	C_{80} (dB)	15.36	0.63	1	17.18	0.38	1.79	16.42	0.22	1.06
Reverberation Time	T_{30} (s)	0.20	0.012	5% = 0.01	0.15	0.0025	0.05	0.20	0.001	0.0
Gravity Time	T_s (ms)	20.69	1.86	10	19.01	1.02	1.89	20.08	0.57	0.61

Reconstructed RIR preserves the perceptual characteristics \Leftarrow Deviations ($e_{R,M}$) are within the Just Noticeable Difference (JND) limits

Current Work

Testing for different rooms

- Room of size $6.5 \times 8.3 \times 2.9$ m and $T_{60} = 1.12$ s
- $\mathbf{y}_o = (1, 90^\circ, 0^\circ)$
- Reconstructed for a receiver at $\mathbf{x} = (0.042, 69^\circ, 0^\circ)$
 - Maximum $A(t, k)$ set to 10
 - Total of 18129 reflection sources $\mathbf{z}_a(t, k)$
- Room of size $5.75 \times 7.87 \times 2.91$ m and $T_{60} = 1.2$ s
- $\mathbf{y}_o = (1.8, 88^\circ, 56^\circ)$
- Reconstructed for a receiver at $\mathbf{x} = (1.29, 87^\circ, 0^\circ)$
 - Maximum $A(t, k)$ set to 5
 - Total of 9465 reflection sources $\mathbf{z}_a(t, k)$

Parameters	Measured		Reconstructed	
	Value	JND	Value	$e_{R,M}$
EDT (s)	0.36	5% = 0.018	0.35	0.01
C_{80} (dB)	12.57	1	12.23	0.34
T_{30} (s)	0.47	5% = 0.0235	0.47	0.0
T_{ζ} (ms)	16.17	10	17.56	1.39

Parameters	Measured		Reconstructed	
	Value	JND	Value	$e_{R,M}$
EDT (s)	0.35	5% = 0.0175	0.38	0.03
C_{80} (dB)	13.55	1	13.36	0.19
T_{30} (s)	0.33	5% = 0.0165	0.32	0.01
T_{ζ} (ms)	16.34	10	17.53	1.19

Conclusion

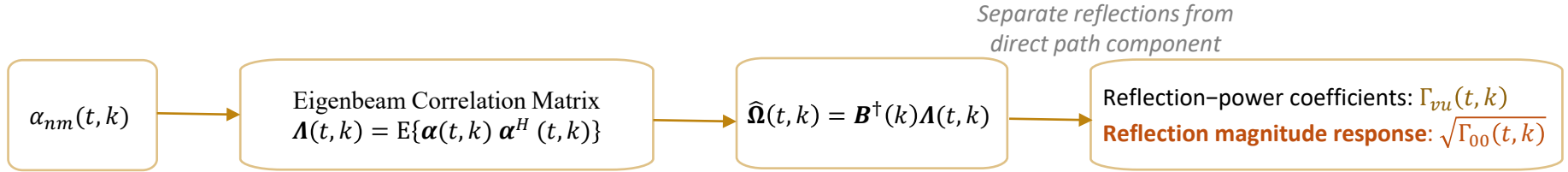
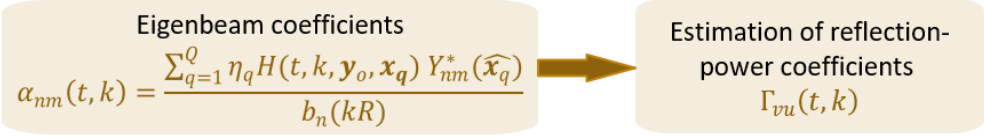
- *Reconstructed RIRs successfully preserved the temporal and spectral behaviors of the measured RIRs*
 - Achieved with **single SMA measurement** and **few parameters per time-frequency frame**
 - Dominant directions of reflections and their relative weights
 - Reflection magnitude response
 - Limitation: Approximate dimensions of the room should be known
- *Future Works:*
 - Incorporating angular spread (κ_a) of reflections from room surfaces
 - Develop the method to facilitate listener translations
 - Perceptual evaluation

Thank you for your attention !

Questions?



Australian
National
University



$$\mathbf{B}(k) = \begin{bmatrix} \delta_{0000} & d_{000000} & \cdots & d_{0000VV} \\ \delta_{001-1} & d_{001-100} & \cdots & d_{001-1VV} \\ \vdots & \vdots & \vdots & \vdots \\ \delta_{00NN} & d_{00NN00} & \cdots & d_{00NNVV} \\ \delta_{1-100} & d_{1-10000} & \cdots & d_{1-100VV} \\ \vdots & \vdots & \vdots & \vdots \\ \delta_{NNNN} & d_{NNNN00} & \cdots & d_{NNNNVV} \end{bmatrix}$$

$$\Lambda(t, k) = \begin{bmatrix} \Lambda_{0000} \\ \Lambda_{001-1} \\ \vdots \\ \Lambda_{00NN} \\ \Lambda_{1-100} \\ \vdots \\ \Lambda_{NNNN} \end{bmatrix}$$

$$\mathbf{\Omega}(t, k) = \begin{bmatrix} P_D \\ \Gamma_{00} \\ \Gamma_{1-1} \\ \vdots \\ \Gamma_{V-V} \\ \vdots \\ \Gamma_{VV} \end{bmatrix}$$

Direct Path Power
 Reflection power response $\propto |H_r(t, k, \mathbf{y}_o, \mathbf{x})|$
 (Direction-independent)
 $\Gamma_{vu}(k, t)$: Reflection Power Coefficients

$$\delta_{nmn'm'} : 16\pi^2 i^{(n-n')} Y_{nm}^*(\widehat{\mathcal{Y}}_o) Y_{n'm'}(\widehat{\mathcal{Y}}_o)$$

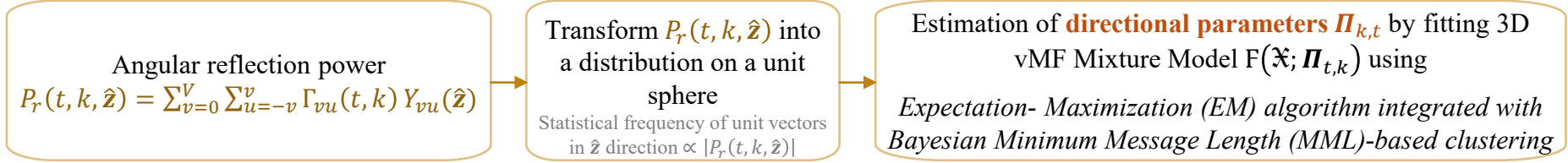
$$d_{nmn'm'vu} : 16\pi^2 i^{(n-n')} (-1)^m \sqrt{\frac{(2v+1)(2n+1)(2n'+1)}{4\pi}} \begin{pmatrix} v & n & n' \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} v & n & n' \\ u & -m & m' \end{pmatrix}$$

$$\Lambda_{nmn'm'} : E\{\alpha_{nm}(t, k) \alpha_{n'm'}^*(t, k)\}$$

Angular reflection power
 $P_r(t, k, \hat{\mathbf{z}}) = \sum_{v=0}^V \sum_{u=-v}^v \Gamma_{vu}(t, k) Y_{vu}(\hat{\mathbf{z}})$

Estimation of directional parameters $\boldsymbol{\mu}_a(t, k)$ and $\omega_a(t, k)$ by fitting 3D von Mises-Fisher (vMF) Mixture Model

$$F(\mathfrak{X}; \boldsymbol{\Pi}_{t,k}) = \sum_{a=1}^{A(t,k)} \omega_a(t, k) \frac{\sqrt{\kappa_a(t,k)}}{(2\pi)^{\frac{3}{2}} I_{\frac{1}{2}}(\kappa_a(t,k))} e^{(\kappa_a(t,k) \boldsymbol{\mu}_a(t,k)^T \mathfrak{X})}$$



$$F(\mathfrak{X}; \boldsymbol{\Pi}_{t,k}) = \sum_{a=1}^{A(t,k)} \omega_a(t, k) \frac{\sqrt{\kappa_a(t,k)}}{(2\pi)^{\frac{3}{2}} I_{\frac{1}{2}}(\kappa_a(t,k))} \exp(\kappa_a(t, k) \boldsymbol{\mu}_a(t, k)^T \mathfrak{X})$$

- \mathfrak{X} : Set of uniformly sampled points on the surface of a unit sphere
- $\boldsymbol{\mu}_a(t, k)$: Mean Direction Vector \Rightarrow *Dominant Reflection Directions* of $\{t, k\}^{th}$ bin
- $\omega_a(t, k)$: Convex mixing coefficients \Rightarrow *Relative strength of each $\boldsymbol{\mu}_a$ direction*
- $A(t, k)$: Number of vMF components \Rightarrow *Number of dominant reflection directions*
- $\kappa_a(t, k)$: Non-negative Concentration Parameter \Rightarrow *Dispersion of reflected power from $\boldsymbol{\mu}_a$*
- $\boldsymbol{\Pi}_{t,k}$: $\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_A; \omega_1, \dots, \omega_A; \kappa_1, \dots, \kappa_A\} \Rightarrow$ *Detectable only for $\{t, k\}$ bins with anisotropic (early) reflections*

