

Diff-SAGE: End-to-End Spatial Audio Generation using Diffusion Models



Saksham Singh Kushwaha^{1,2,*}, Jianbo Ma², Mark R. P. Thomas², Yapeng Tian¹, Avery Bruni²

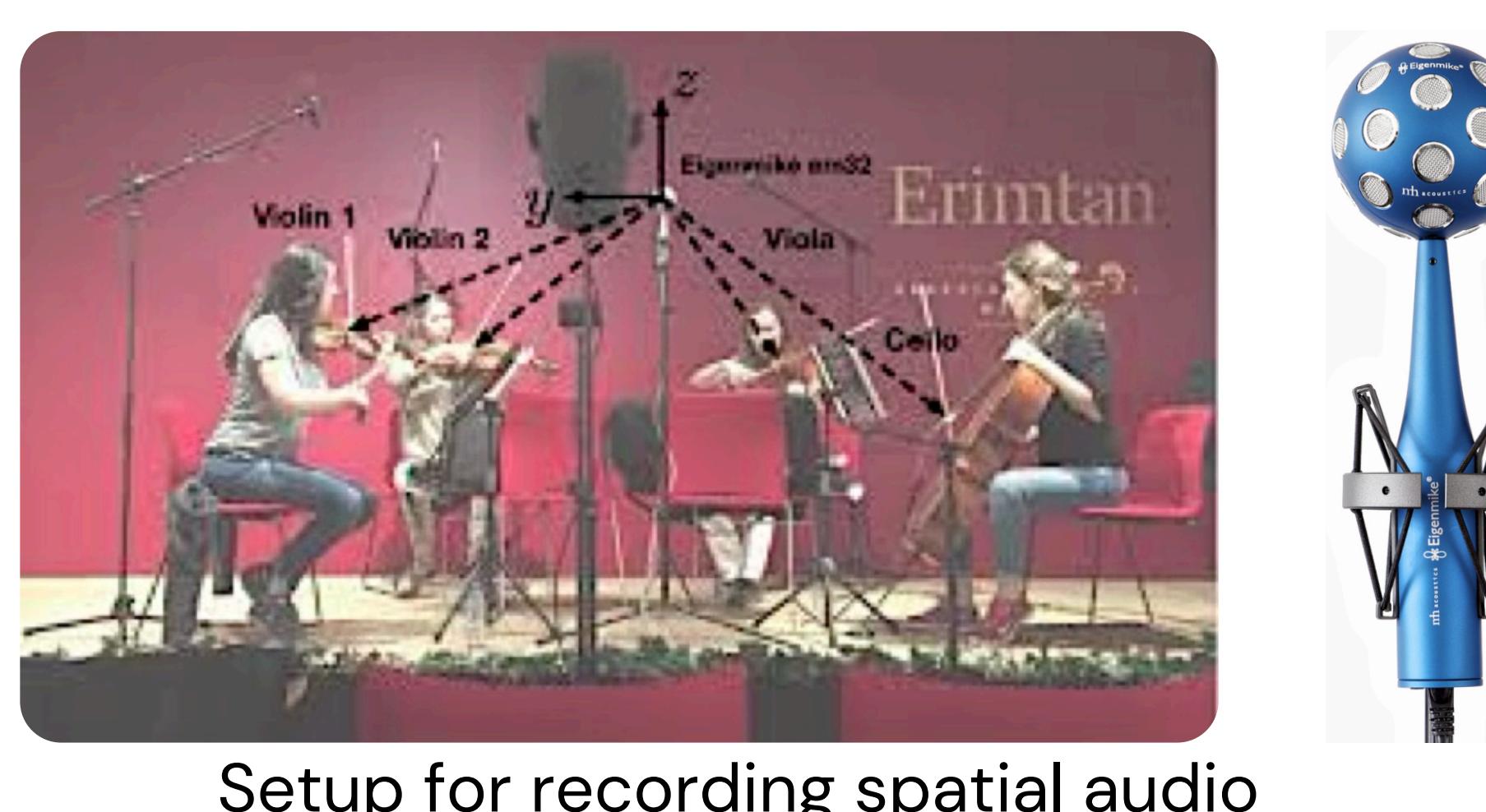
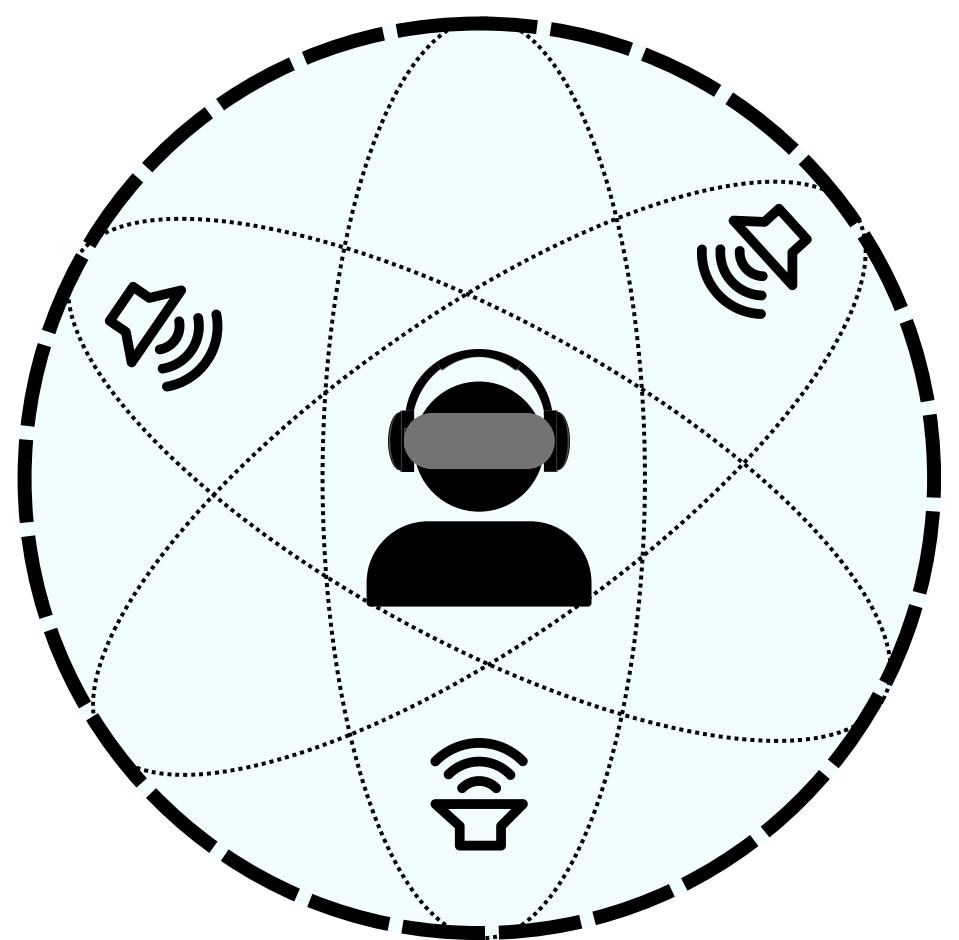
¹University of Texas at Dallas, ²Dolby Laboratories

(*Work done during internship at Dolby)



Introduction

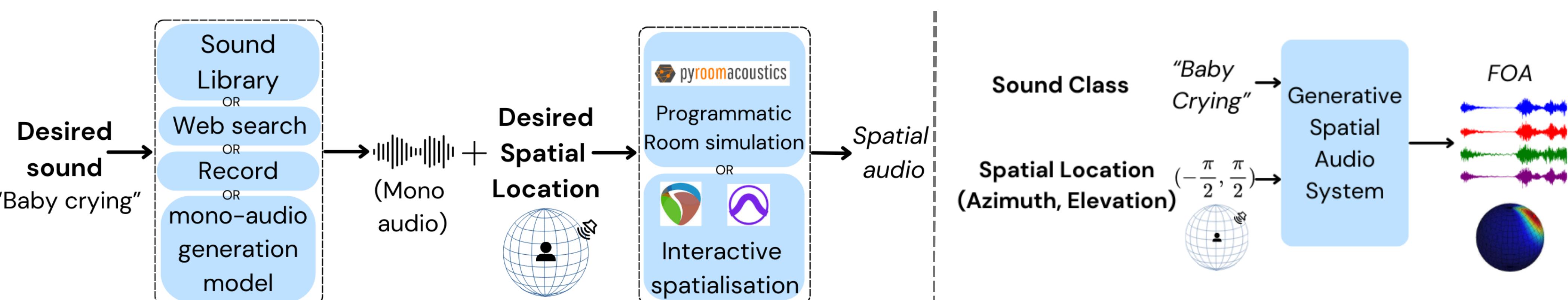
- Spatial audio creates immersive experience.
- Recording is expensive and complex.



Setup for recording spatial audio

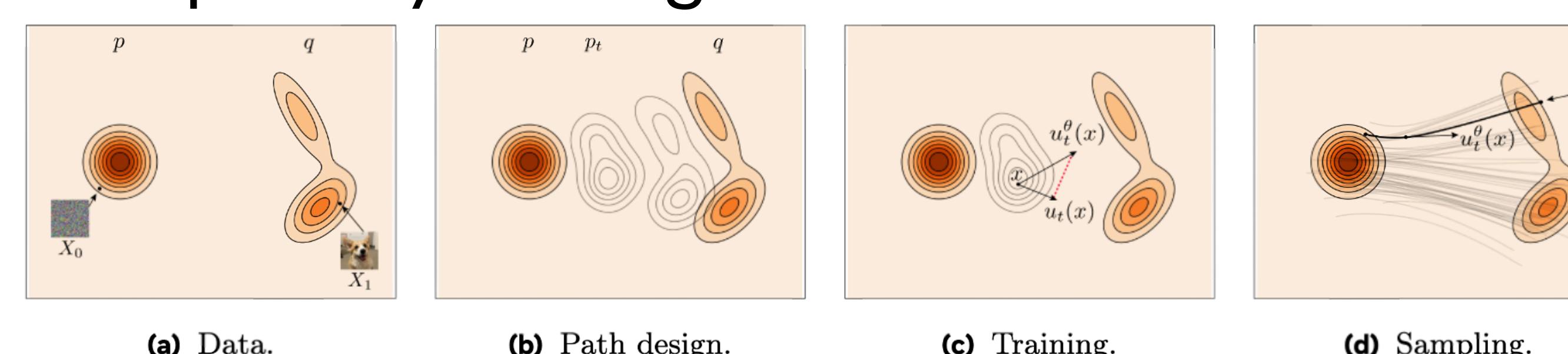
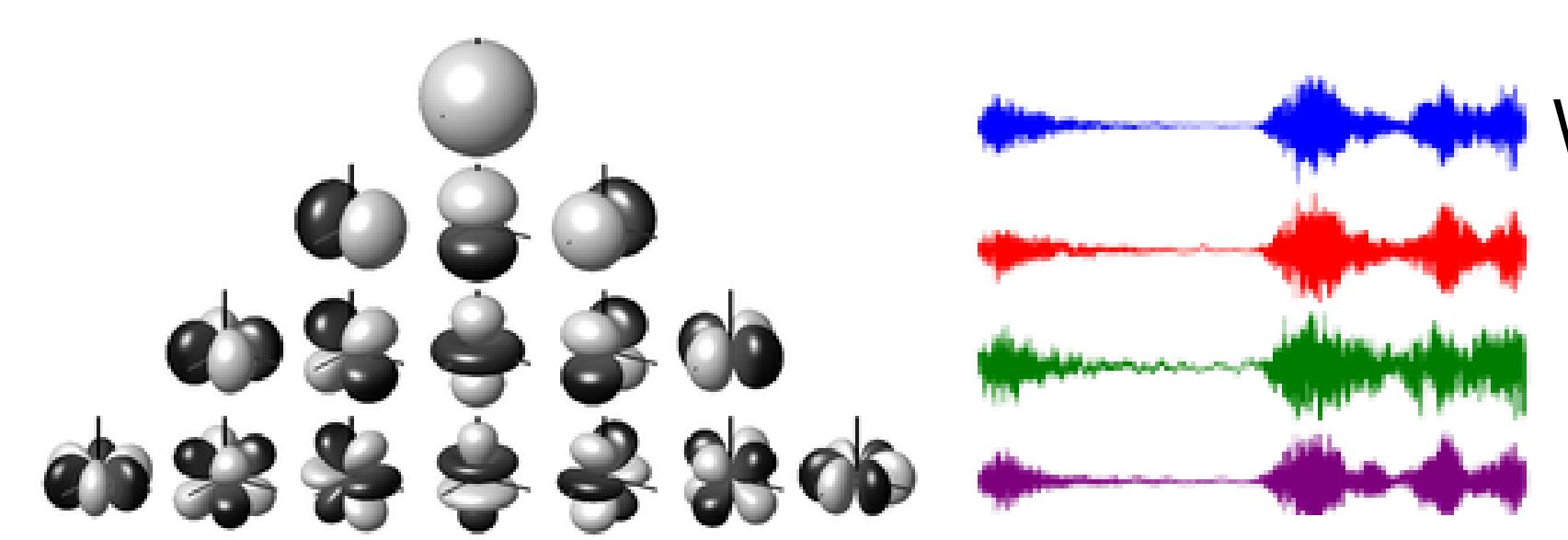
Motivation

- Traditional approaches are not scalable, require expertise etc.
- Our task: end-to-end spatial audio generation
 - Input: Sound class + Spatial location
 - Output: First Order Ambisonics (FOA)



Background

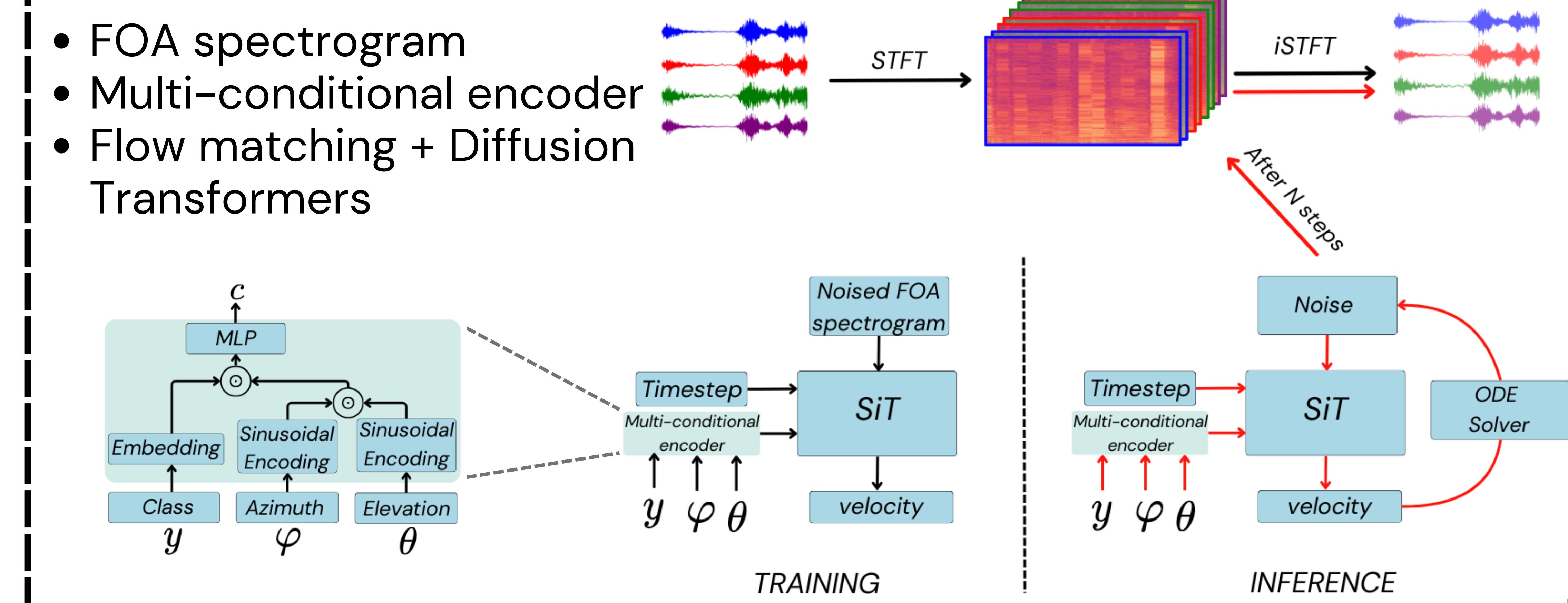
- First Order Ambisonics
 - Mono/Omnidirectional mic: W
 - Directional mics: X, Y, Z
- Flow Matching
 - Learn velocity to map noisy to target data



Method

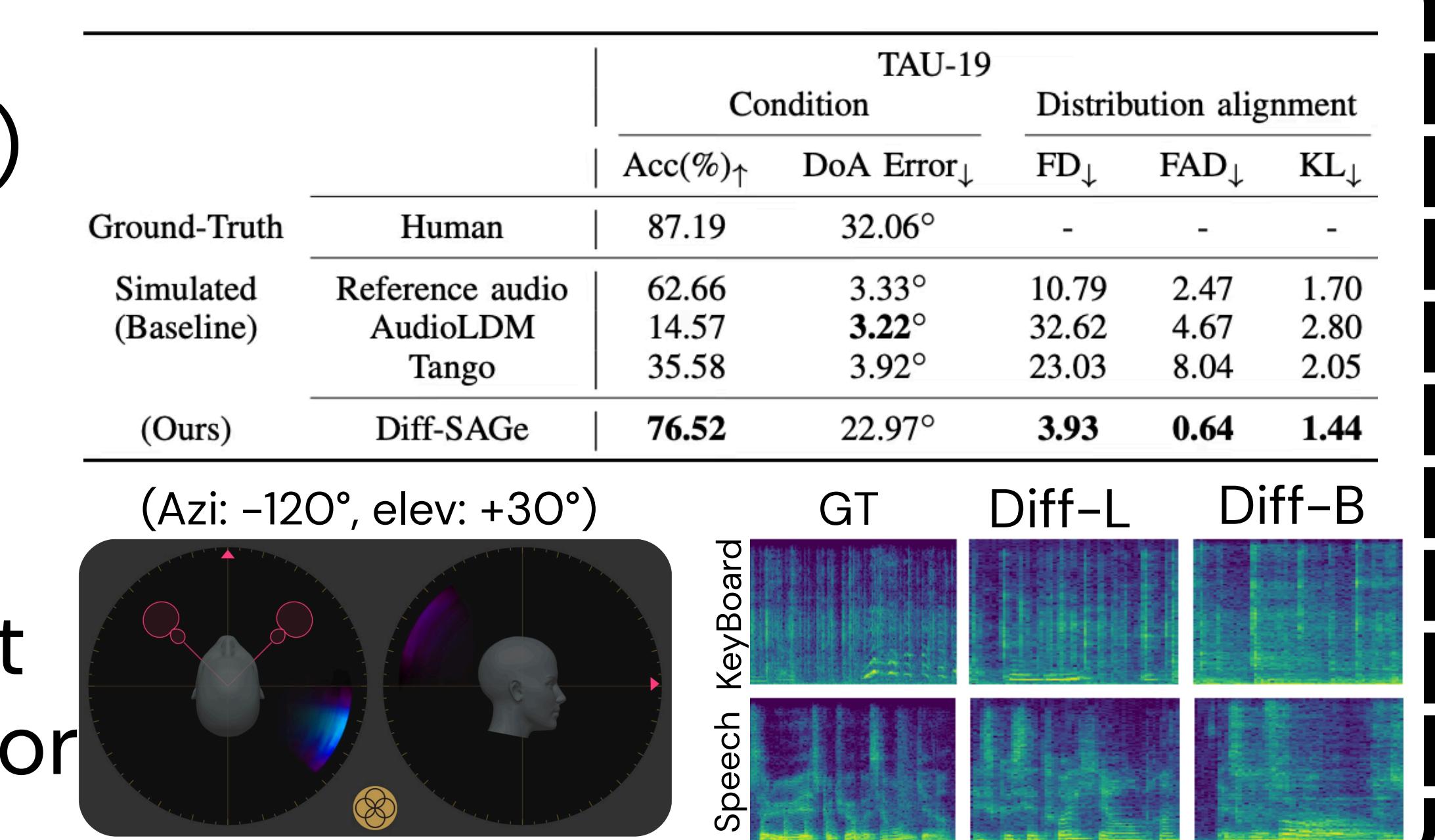
- FOA spectrogram
- Multi-conditional encoder
- Flow matching + Diffusion

Transformers



Results

- Datasets:
 - Sound event localisation (SELD)
 - TAU-19, TAU-20
- Baselines:
 - Mono: Real & Generated
 - Simulation: Pyroomacoustics
- Metrics: Condition + Dist. alignment
- DoA Error → Human annotation error



Conclusion

- Introduce a new problem of spatial audio generation.
- Proposed Diffusion transformer (SiT) w/ multi-conditional encoder.
- Surpass simulation based approaches.
- Future directions:
 - Multiple/moving sources; longer duration; 360° video → FOA

[1] Scheibler et al., "Pyroomacoustics: A Python package for audio room simulation and array processing algorithms," ICASSP 2018.
 [2] Zotter et al., "Ambisonics: A practical 3D audio theory for recording, studio production, sound reinforcement, and virtual reality," Springer Nature, 2019.
 [3] Ma et al., "SiT: Exploring flow and diffusion-based generative models with scalable interpolant transformers," ECCV 2024.
 [4] Lipman et al., "Flow Matching Guide and Code", arXiv 2024.
 [5] Advanee et al. "A multi-room reverberant dataset for sound event localization and detection", arXiv 2019
 [6] DAW and plugin credits: "Ambisonics audio processor Soundfield by Rode" and "Reaper"