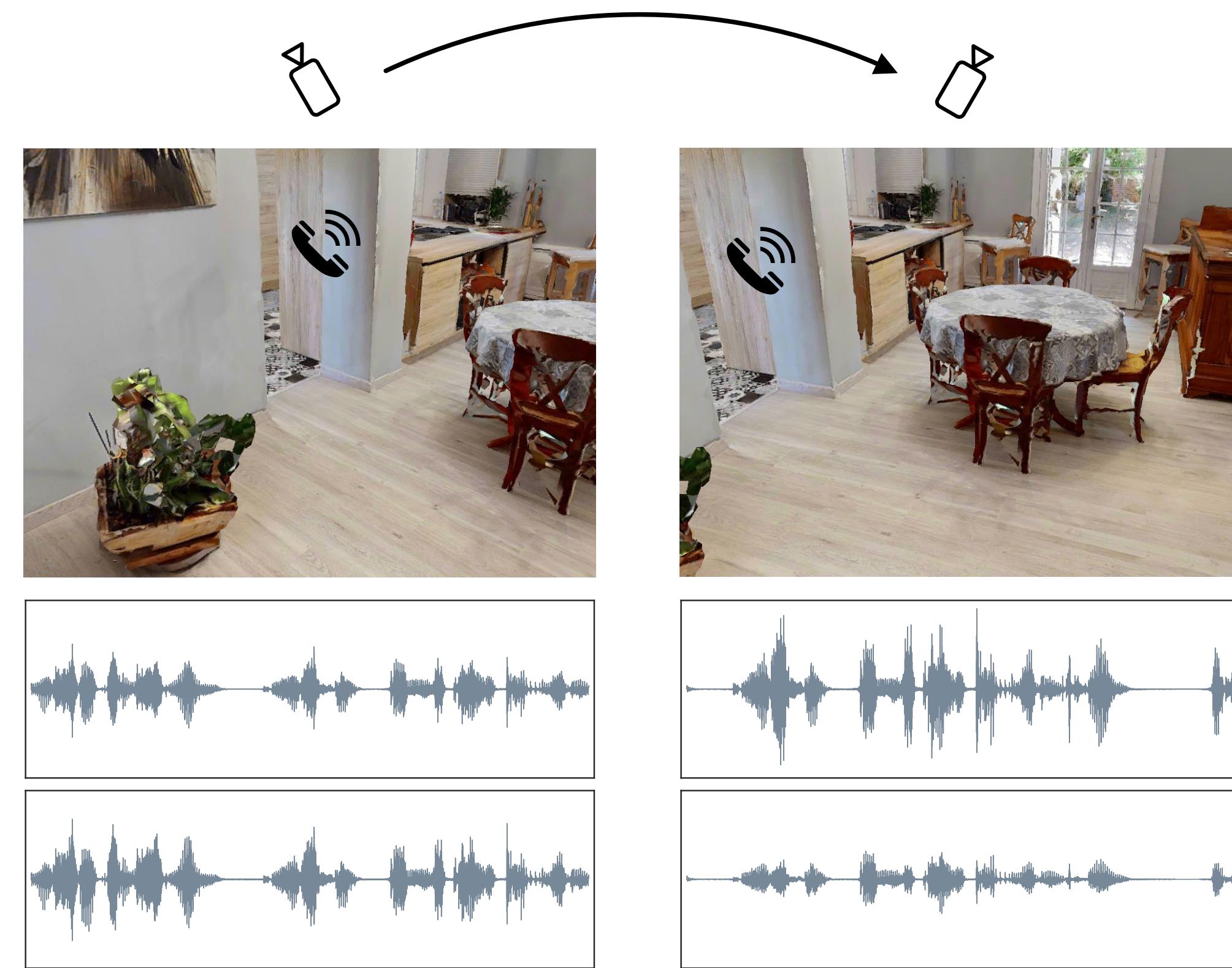


Sound Localization from Motion: Jointly Learning Sound Direction and Camera Rotation

Ziyang Chen, Shengyi Qian, Andrew Owens
 University of Michigan

Motivation

The images and sounds that we perceive undergo subtle but geometrically consistent changes as we rotate our heads. Can we use these cues to learn audio and visual models of space?

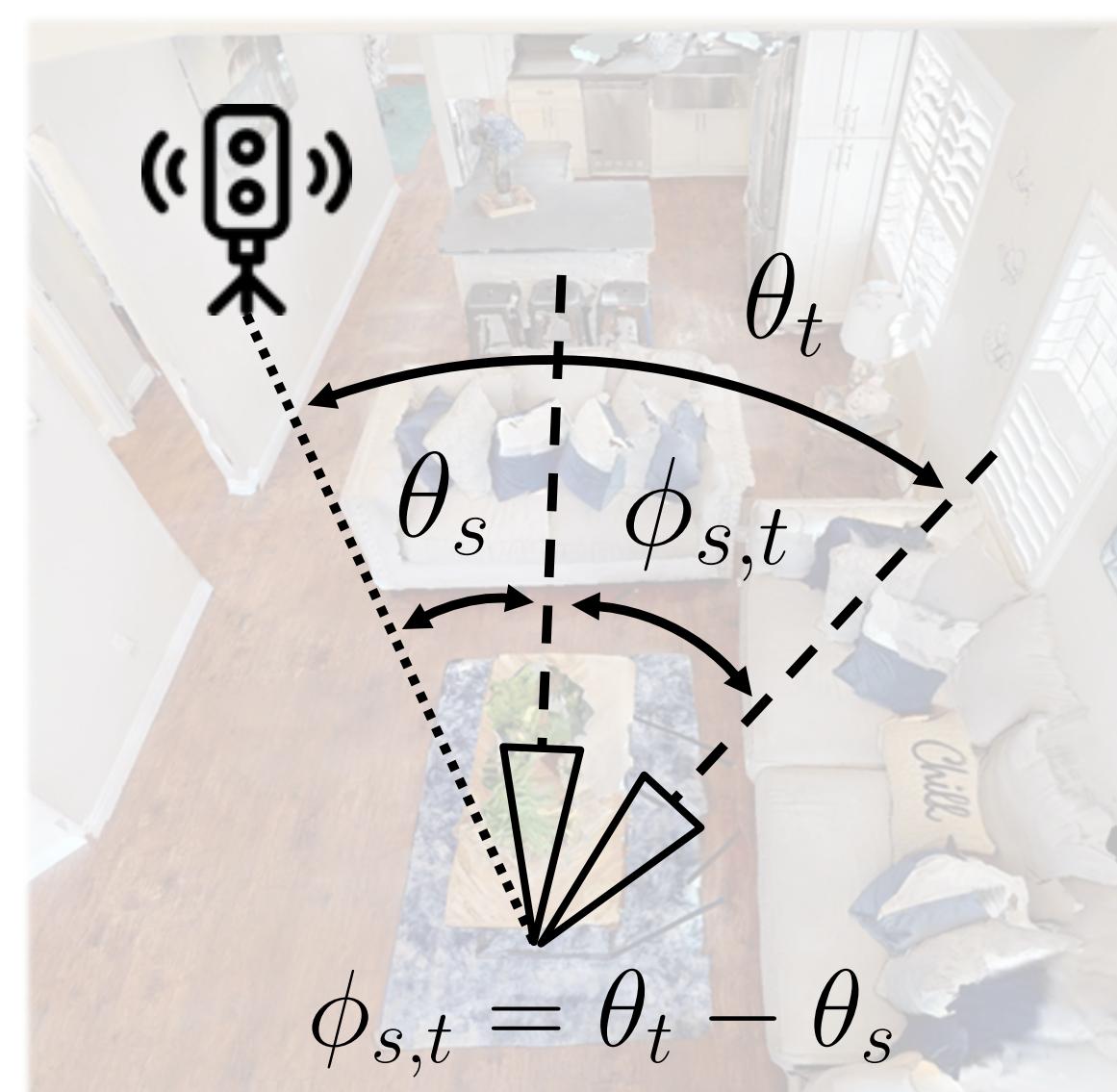


Sound Localization from Motion

We propose SLfM: jointly learning sound direction and camera rotation from multi-view audio-visual data.

Idea: learn to enforce cross-modal consistency. Audio and visual predictions should agree with each other:

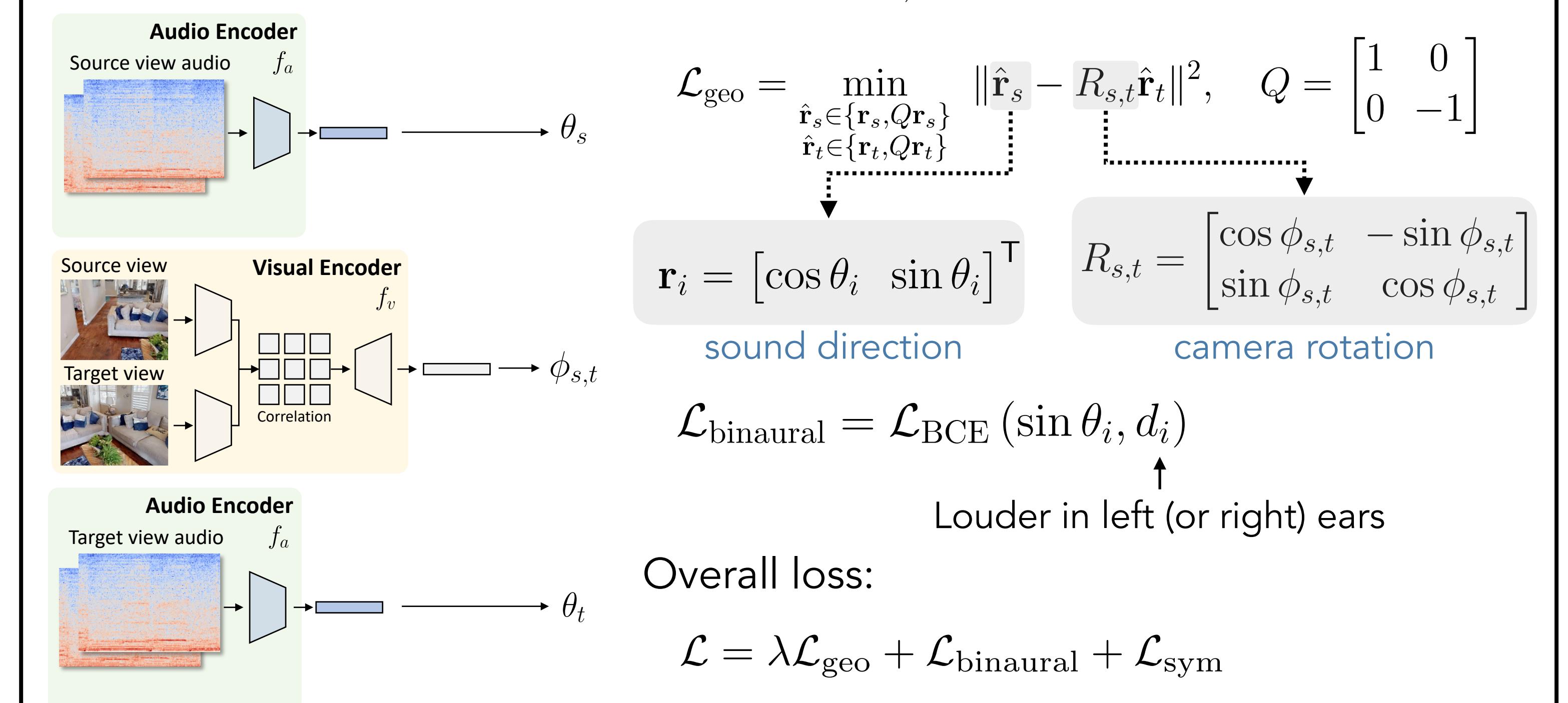
- **Audio model:** predict sound direction from stereo audio.
- **Visual model:** predict camera rotation from two images.



Method

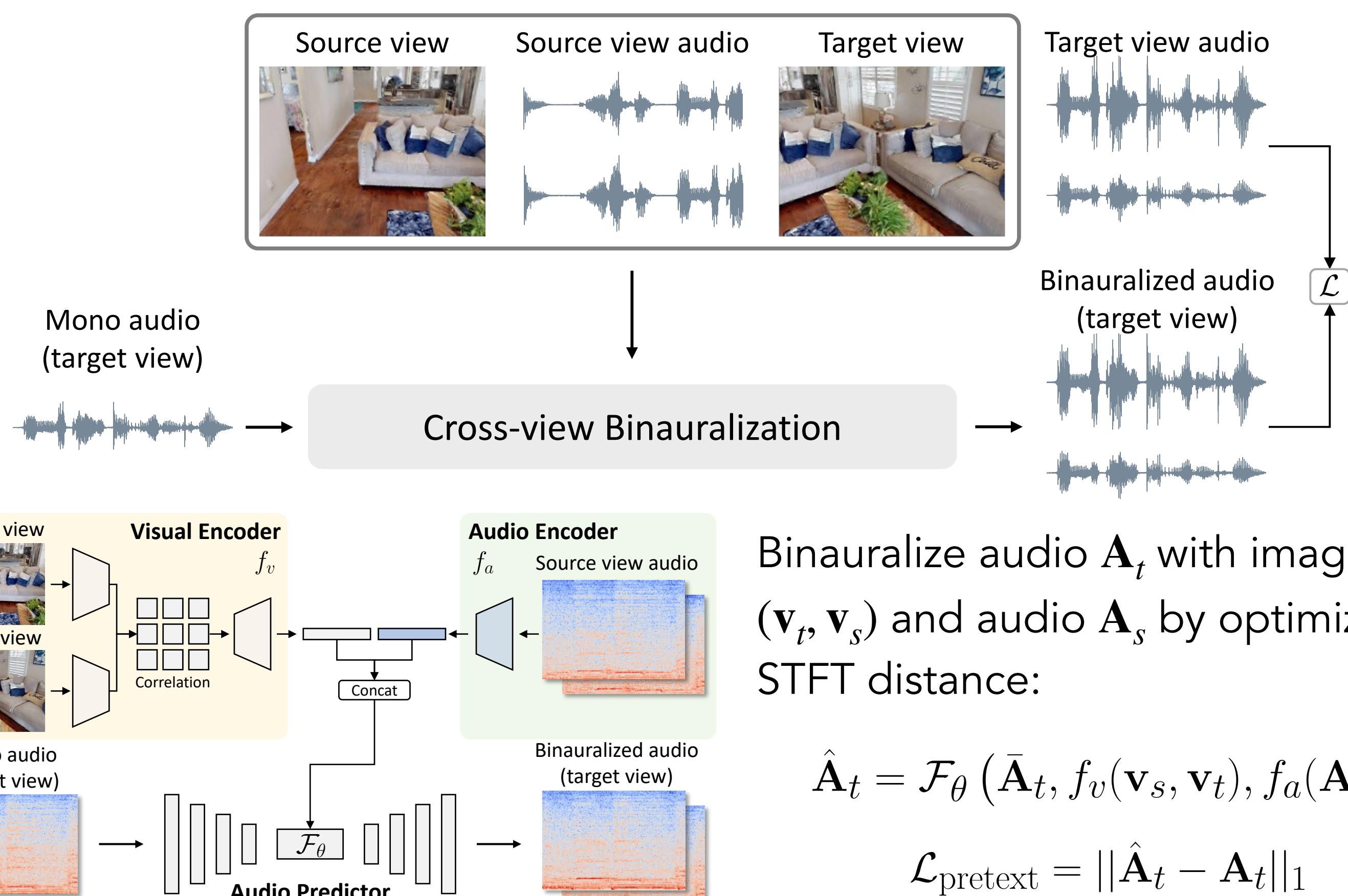
Estimating pose and localizing sound

Cross-modal geometric consistency: $\phi_{s,t} = \theta_t - \theta_s$



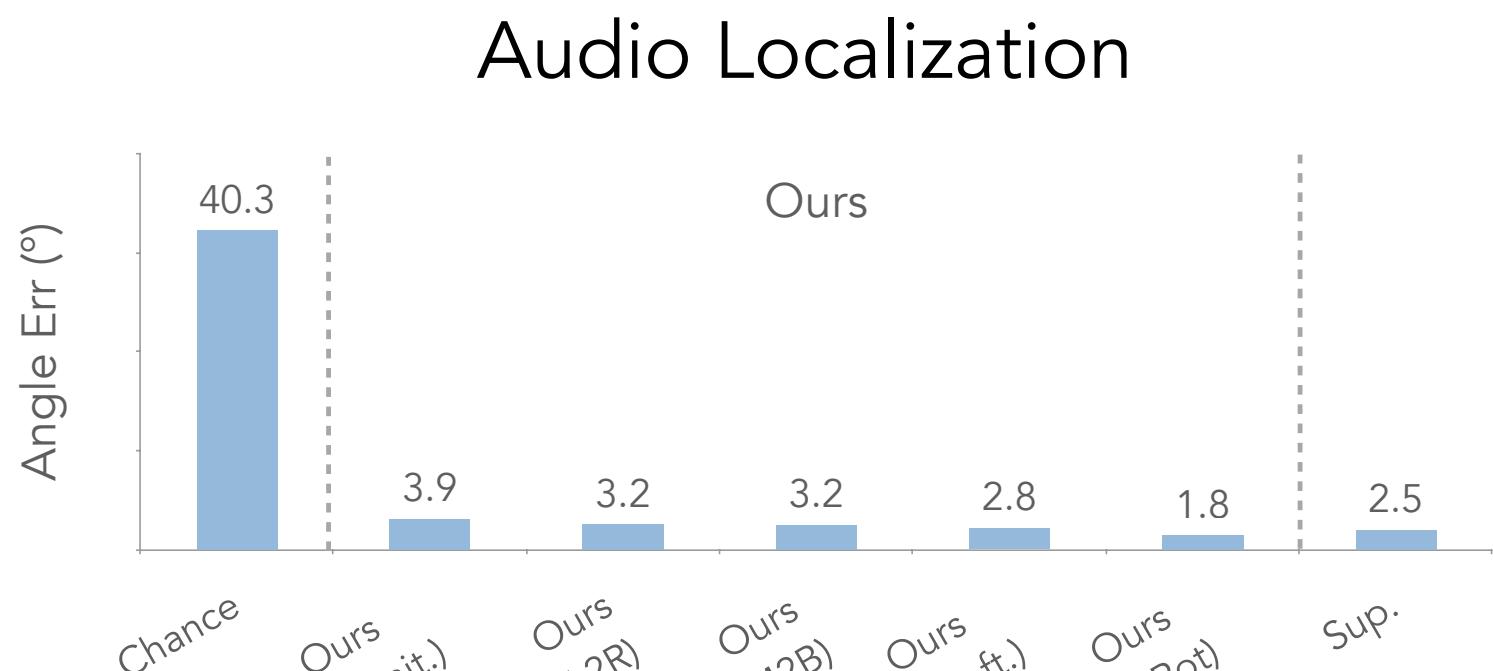
Learning representation via spatialization

We learn an audio-visual representation that conveys spatial cues by solving a cross-view binauralization task.

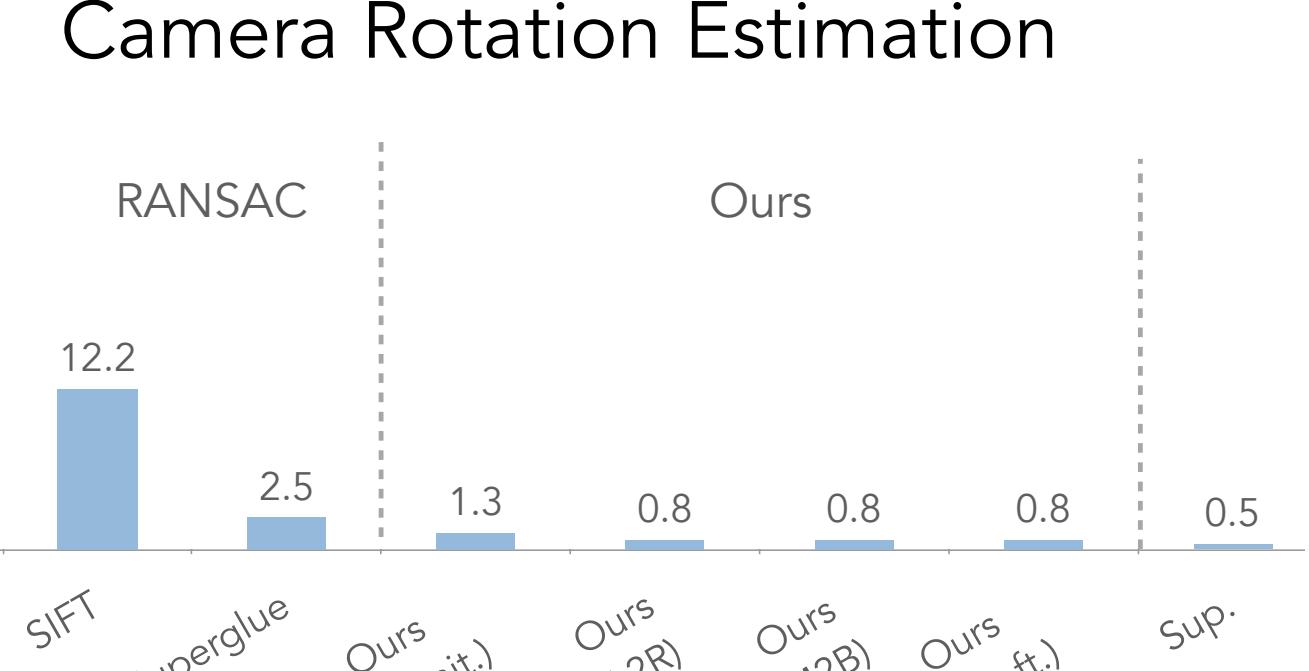


Experiments

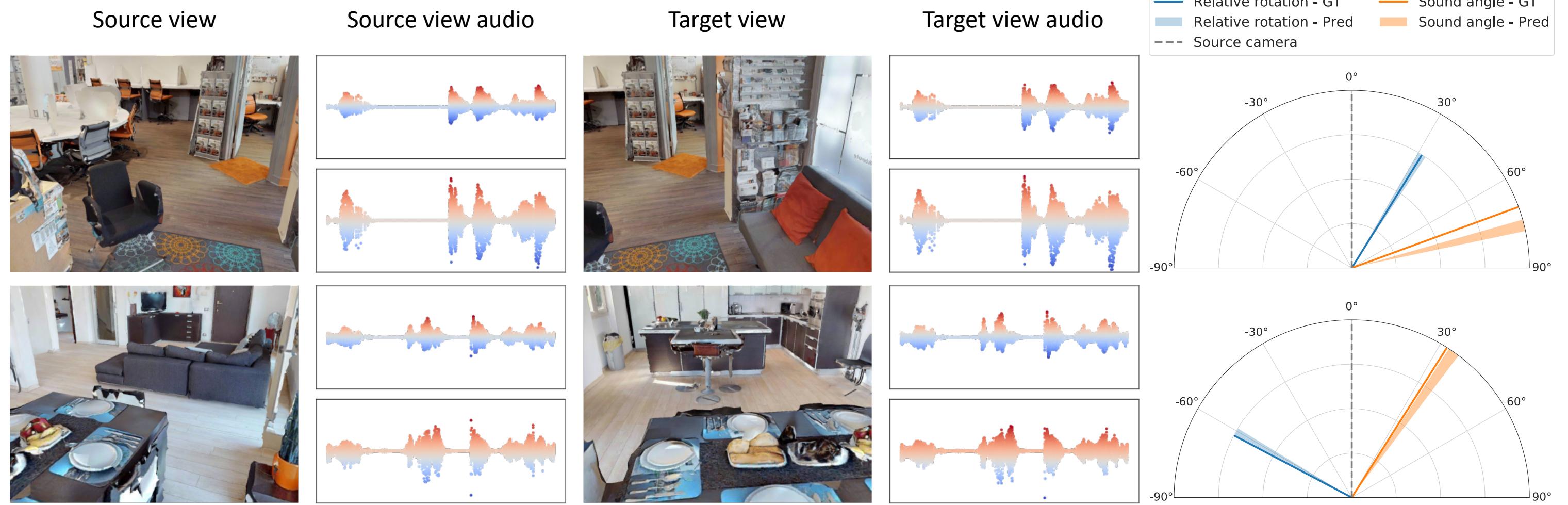
Audio Localization



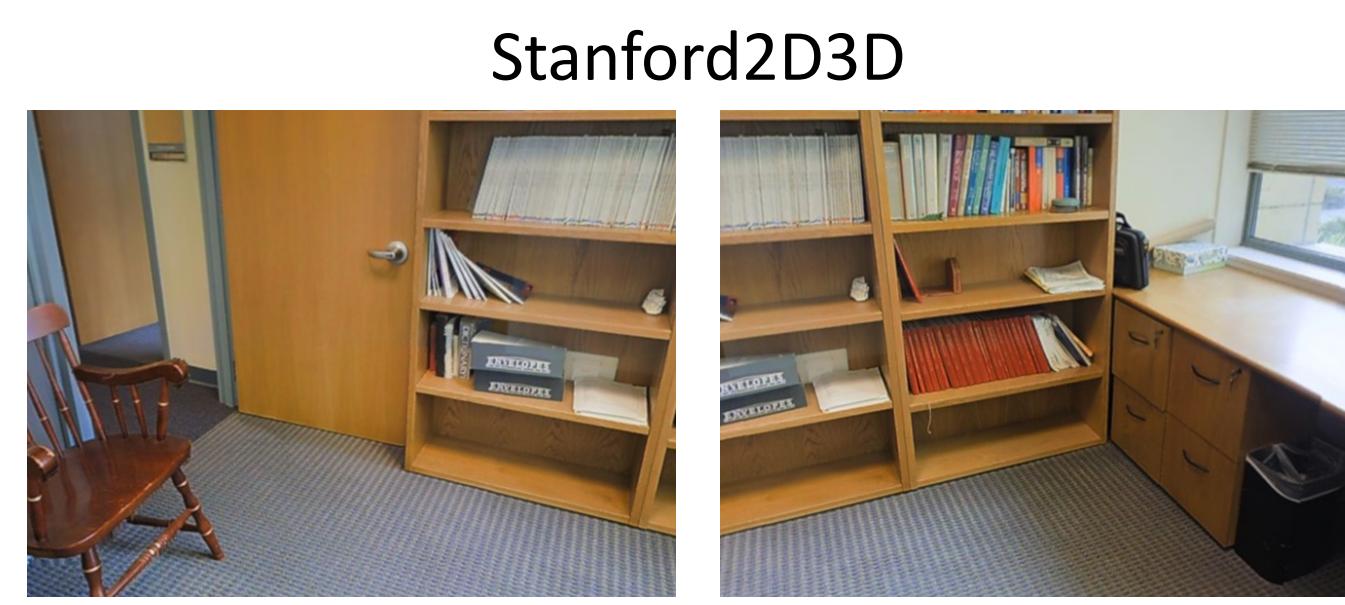
Camera Rotation Estimation



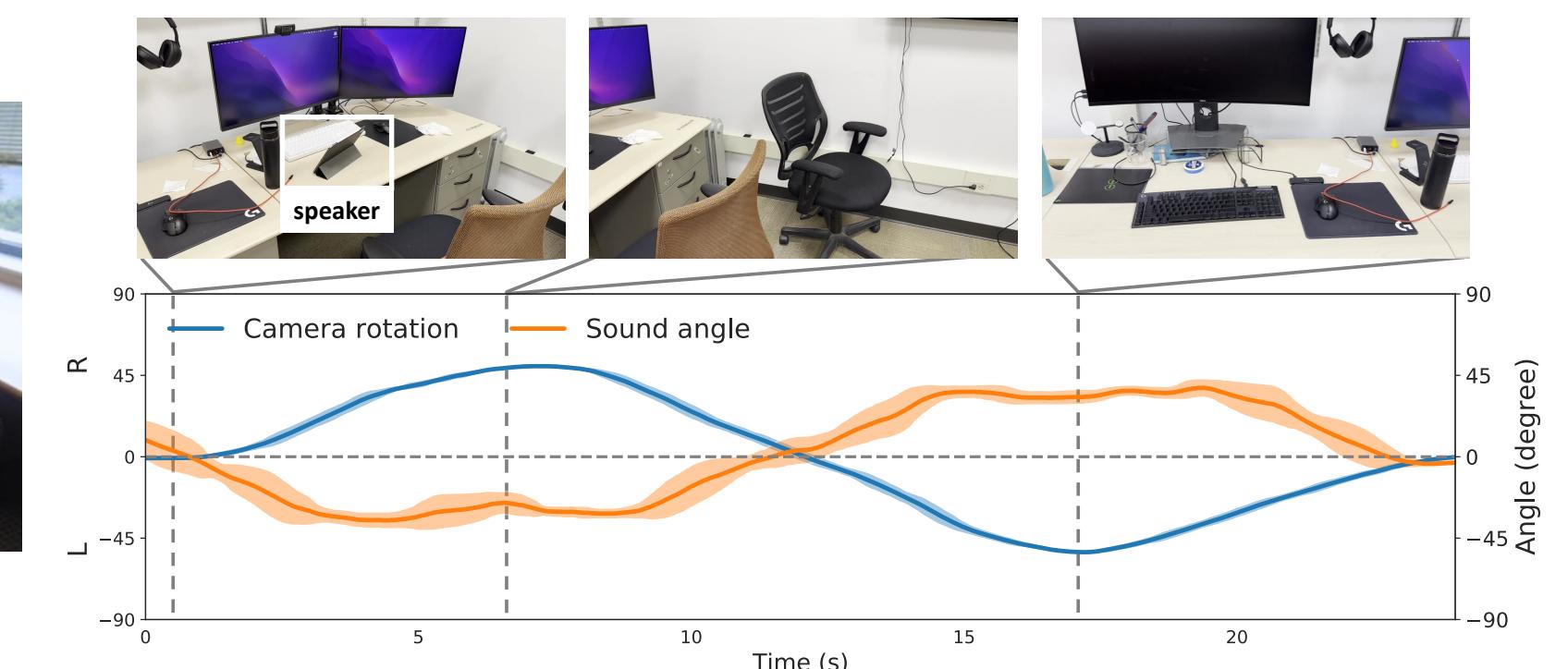
Qualitative results



In-the-wild examples



Real-world demo



Related work

- Yang, Karren, et al. "Camera Pose Estimation and Localization with Active Audio Sensing." ECCV 2022
- Chen, Changan, et al. "Novel-View Acoustic Synthesis." CVPR 2023.
- Chen, Changan, et al. "Soundspace 2.0: A simulation platform for visual-acoustic learning." NeurIPS 2022