Audio-Visual Embodied AI: From Simulating to Navigating with Sounds in Spaces

Changan Chen changan.io UT Austin 10/03/2023



Human perception is multisensory

We often use vision, audio, touch, smell to sense the world





Autonomous agents

Home assistance robot



Rescue robot



Robots that can **see, hear and act** in the environment, e.g., reacting to humans' commands or locating fire-alarm going off

Simulating embodiment in 3D scenes

Datasets







Savva et al., ICCV 2019



Xia et al., ICRA 2020



Kovle et al., arXiV 2017

Advantages: Large-scale training, fast experimentation, consistent benchmarking and replicable research

Sim2Real

Sim2Real Predictivity: Does Evaluation in Simulation Predict Real-World Performance, Kadian et al., IRAL 2020 Sim-to-Real Transfer for Vision-and-Language Navigation, Anderson et al., CoRL 2020 RoboThor: An Open Simulation-to-Real Embodied AI Platform, Deitke et al., CVPR 2020

Enabling embodied agents and tasks



Today's embodied agents (robots) are deaf

• We want robots that can hear and react in the environment

- No existing simulation supports audio-visual rendering
- No existing formulation for audio-visual navigation

SoundSpaces demo

Background: acoustic simulation

Goal: simulate a perceptuallyvalid approximation of the room impulse response (RIR)

↑

Energy

Physics-based audio rendering

3D Geometry

Material Properties

Simulate the sound received by the listener from a source location

Sound propagation system

3D spatial audio for reflections and reverb with realistic acoustics based on bidirectional ray tracing

Real-scan environments

Replica¹ dataset

Matterport3D² dataset

¹The Replica Dataset: A Digital Replica of Indoor Spaces, Straub et al., arXiv, 2019 ²Matterport3D: Learning from RGB-D Data in Indoor Environments, Chang et al., 3DV, 2017

SoundSpaces: our audio simulator

SoundSpaces produces realistic audio rendering based on the room geometry, materials, and sound source location by precomputing the room impulse response function (RIR)

Users can insert any sound of their choice at runtime. The received sound is obtained by convolving the RIR with the source sound.

	# Scenes	Avg. Area	# RIRs
Replica	18	47.24 m ²	0.9M
Matterport3D	85	517.34 m ²	16.7M

Table: Summary of dataset statistics

Visit soundspaces.org for more information!

Enabling audio-visual embodied AI and beyond

(e)

(f)

SoundSpaces 2.0: A fast, continuous, configurable and generalizable audio-visual simulation platform

Continuous rendering

We offer both spatial and acoustic continuity.

Navigating to someone speaking

Changan Chen et al., SoundSpaces 2.0: A Simulation Platform for Visual-Acoustic Learning, NeurIPS 2022

Configurable simulation

Users can change all these parameters!

Simulation parameters

- Frequency bands
- Direct sound
- Indirect sound
- Transmission
- Diffraction
- Number of rays
- Number of threads
- Sample rate

Microphone types

- Mono
- Binaural
- Stereo
- Quad
- Surround_5_1
- Surround_7_1
- Ambisonics
- Your mic array
- ...

Material properties

- Absorption coefficients
- Scattering coefficients
- Transmission coefficients
- Damping coefficients
- Frequency band specs
- Instance level config
- •

Generalizable simulation

We support arbitrary scene datasets.

Changan Chen et al., SoundSpaces 2.0: A Simulation Platform for Visual-Acoustic Learning, NeurIPS 2022

Validating simulation with real IRs

We collect acoustic measurements of the apartment in Replica dataset and compare to IRs rendered in SoundSpaces

SoundSpaces 2.0 has a better match of direct-to-reverberant ratio with real

Changan Chen et al., SoundSpaces 2.0: A Simulation Platform for Visual-Acoustic Learning, NeurIPS 2022

Main differences

SoundSpaces 1.0

- 500 fps+
- Discrete and unconfigurable

SoundSpaces 2.0

- 30 fps+
- Continuous and configurable

Audio-visual navigation in 3D environments

An agent navigates to a sounding object with vision and audio perception

Learning with deep reinforcement learning

- Learn to navigate in simulation via trials and errors
- Rewarded +1 for getting close and +10 for reaching the goal

Navigation policy

Navigation example

The agent leverages the complementary spatial information in audio and vision, and navigates to the goal successfully

Limitations of the navigation policy

Existing models learn to act at fixed granularities of action motion

- Chen et al.¹: learn to generate primitive actions step-by-step
- Gan et al.²: predict target locations and navigate with geometric planner

¹SoundSpaces: Audio-Visual Navigation in 3D Environments, Chen et al., ECCV, 2020 ²Look, Listen, and Act: Towards Audio-Visual Embodied, Gan et al., ICRA, 2020

Learning to set waypoints for AV navigation

- Infer audio-visual subgoals with RL end-to-end at varying granularities
- Acoustic memory to help infer goal locations and decide stop actions

Audio-visual waypoints navigation model (AV-WAN)

Changan Chen et al., Learning to Set Waypoints for Audio-Visual Navigation, ICLR 2021

Waypoint selection and acoustic memory

Changan Chen et al., Learning to Set Waypoints for Audio-Visual Navigation, ICLR 2021

Navigation results

- Our model (AV-WaN) strongly outperforms all baselines and existing methods
- · Generalizing to unheard sounds is much more challenging

Changan Chen et al., Learning to Set Waypoints for Audio-Visual Navigation, ICLR 2021

Navigation trajectories

- Gan et al. [ICRA 20]: is prone to errors and often leads the agent to backtrack
- Chen et al. [ECCV20]: oscillates around obstacles
- AV-WaN (Ours): reaches the goal most efficiently

Changan Chen et al., Learning to Set Waypoints for Audio-Visual Navigation, ICLR 2021

Robustness to audio noise

AV-WaN with acoustic map is more robust to audio noise

Changan Chen et al., Learning to Set Waypoints for Audio-Visual Navigation, ICLR 2021

Limitations of the AudioGoal task

AudioGoal task (Chen et al. ECCV 2020, Gan et al. ICRA 2020):

- The sound is constant and periodic (it covers the whole episode)
- The goal has no visual embodiment

The agent searches for the ringing telephone in an unfamiliar environment

The agent must continue navigating even after the sound stops

Our proposed semantic AudioGoal task:

- The sound is associated with a semantically meaningful object
- The sound is not periodic and has variable length

Semantic AudioGoal dataset

- Augment an existing simulator SoundSpaces¹ with semantic sounds
- 21 object categories in Matterport3D²: chair, TV, cabinet, sink etc.
- Object-emitted sounds and object-related sounds

¹Changan Chen et al., SoundSpaces: Audio-Visual Navigation in 3D Environments, ECCV 2020 ²Angle Chang et al., Matterport3D: Learning from RGB-D Data in Indoor Environments, 3DV 2017

Semantic audio-visual navigation

- Learn the association between how objects look and how they sound
- Leverage long-term memory to handle sporadic acoustic events

Changan Chen et al., Semantic Audio-Visual Navigation, CVPR 2021

Semantic audio-visual navigation model (SAVi)

Navigation results

- SAVi strongly outperforms all existing methods
- Generalizing to unheard sounds is much more challenging

Changan Chen et al., Semantic Audio-Visual Navigation, CVPR 2021

Navigation example

The agent identifies it's drawer sliding sound and locates the target object with vision after the sound stops.

Collaborative acoustic measurements

Measuring acoustics traditionally relies on manual sampling

Much more to be explored ...

Audio-Visual Navigation Chen et al., ECCV 2020

Echolocation Learning Gao et al., ECCV 2020

Audio-Visual Separation Majumder et al., ICCV 2021

Audio-Visual Mapping Purushwalkam et al., ICCV 2021

Neural Acoustic Rendering Luo et al., NeurIPS 2022

Audio Scene Reconstruction Chen et al., ICCV 2023

