SoundSpaces: Audio-Visual Navigation in 3D Environments

Changan Chen^{*1,4}, Unnat Jain^{*2,4}, Carl Schissler³, Sebastia V. Amengual Gari³, Ziad Al-Halah¹, Vamsi K. Ithapu³, Philip Robinson³, Kristen Grauman^{1,4}

¹UT Austin, ²UIUC, ³Facebook Reality Labs, ⁴Facebook AI Research





Embodied Perception Is a Multisensory Experience

We often use vision, audio, touch, smell to move around Today's agents are deaf!

Vision-Only Gupta et al., 2017 Zhu et al., 2017 Sava et al., 2019		Vision-Language Anderson et al., 2018 Wang et al., 2018 Wang et al., 2019	
	Vision-Interaction Zhu et al., 2017 Gordon et al., 2018 Wortsman et all, 2019	Vision-Audio Chen and Jain et al., 2020 (this work)	

Our contribution: audio-visual embodied navigation --- task and simulation

Audio-Visual Navigation in 3D Environments

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





SoundSpaces: Our Audio Simulator

 We introduce SoundSpaces, an audio simulation platform to enable audio-visual navigation for two visually realistic 3D environments: Replica¹ and Matterport3D²



	# Scenes	Avg. Area	# Training Eps.
Replica	18	47.24 m ²	0.1M
Matterport3D	85	517.34 m ²	2M

Table: Summary of dataset statistics

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

- We introduce SoundSpaces, an audio simulation platform to enable audio-visual navigation for two visually realistic 3D environments: Replica¹ and Matterport3D²
- Our audio simulator produces realistic audio rendering based on the room geometry, materials, and sound source location
- The platform can play varying sounds of your choice in real time by precomputing a transfer function between locations

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



Example 1: Where Is My Phone?



Direction: left ear is louder when the agent faces upward on the top-down map Intensity: overall intensity gets higher as the agent gets closer to the goal

Agent Goal Start Shortest path Agent path Seen/Unseen area Occupied area



Example 2: Where Is The Piano?

Agent view

Top-down map (unknown to the agent)



Audio-Visual Navigation Tasks



AudioGoal



The agent receives an audio signal emitted by the sounding object at each time step

AudioPointGoal

The agent receives both a displacement vector (Δx , Δy) and an audio signal at each time step



Deep RL for Audio-Visual Navigation



Navigation Demo - AudioPointGoal



SPL: 1.00





Navigation Trajectory Comparison







SPL: 0.68

PointGoal agent gets confused about the direction and gets stuck behind the bed.



AudioGoal agent figures out the sound comes from the front more quickly than the PointGoal agent



AudioPointGoal agent knows immediately it should go straight and then right and thus follows the shortest path



Does Audio Help Navigation?

Comparing PointGoal (PG) and AudioPointGoal (APG):

• Audio improves accuracy significantly



Metric: SPL (success weighted by inverse path length)

C. Chen*, U. Jain*, et al., SoundSpaces: Audio-Visual Navigation in 3D Environments, ECCV 2020

12

Can Audio Supplant GPS for AudioGoal?

- AudioGoal is immune to GPS noise (localization error) and robust to microphone noise
- AudioPointGoal degrades less in the presence of GPS noise
- Audio signal gives similar or even better spatial cues than the PointGoal displacements





Effect of Different Sound Sources

From same sound to varied heard sounds to varied unheard sounds¹

- AudioGoal accuracy declines with varied heard sounds to unheard sounds
- AudioPointGoal almost always outperforms AudioGoal agent



¹102 copyright-free sounds, divided into 73/11/18 for train/val/test

C. Chen*, U. Jain*, et al., SoundSpaces: Audio-Visual Navigation in 3D Environments, ECCV 2020

14

What Do the Learned Audio Features Capture?



T-SNE of audio features from an AudioGoal ager

15

Relative Importance of Audio and Vision

Each modality plays an important role in action selection, based on the environment context and goal placement





Conclusion

- Introduce task of audio-visual navigation in 3D environments
- Generalize a state-of-the-art deep RL model
- Introduce SoundSpaces: enabling audio rendering for Habitat
- Create a benchmark suite of tasks for audio-visual navigation

SoundSpaces: Audio-Visual Navigation in 3D Environments

Changan Chen^{*1,4}, Unnat Jain^{*2,4}, Carl Schissler³, Sebastia V. Amengual Gari³, Ziad Al-Halah¹, Vamsi K. Ithapu³, Philip Robinson³, Kristen Grauman^{1,4}

¹UT Austin, ²UIUC, ³Facebook Reality Labs, ⁴Facebook AI Research

Code and audio simulation data available at: <u>http://vision.cs.utexas.edu/projects/audio_visual_navigation</u>





Audio-Visual Waypoints for Navigation

Changan Chen^{1,2}, Sagnik Majumder¹, Ziad Al-Halah¹, Ruohan Gao^{1,2}, Santhosh Kumar Ramakrishnan^{1,2}, Kristen Grauman^{1,2}

UT Austin¹, Facebook AI Research²



Model



Figure 2: Model architecture. Our audio-visual navigation model uses the egocentric stream of depth images and binaural audio (B_t) to learn geometric (G_t) and acoustic (A_t) maps for the 3D environment. The multi-modal cues and partial maps (left) inform the RL policy's prediction of intermediate waypoints (center). For each waypoint, the agent plans the shortest navigable path (right). From this sequence of waypoints, the agent reaches the final AudioGoal efficiently.

Waypoint Selection and Acoustic Memory



Our model dynamically selects waypoints and builds an acoustic memory as it moves.



Start

Waypoint

Normalized intensity

Seen/Unseen area



Experiment Setup

Baselines:

- Random: randomly selects actions and automatically stop when it reaches the goal
- Audio Direction Follower: predicts the audio direction of arrival and moves in that direction
- Frontier Waypoints: intersects the predicted DoA with frontier of the explored area
- Chen and Jain et al.: uses end-to-end RL and predicts primitive actions at each step
- Gan et al.: trains goal predictor and uses geometric planner to reach the goal

Metrics:

- Success (SR): the fraction of successful episodes
- Success weighted by path length (SPL): weighs success with path length
- Normalized distance to goal (NDG): distance to goal relative to the shortest path
- Number of actions (NA): average number of actions per episode
- Success weighted by number of actions (SNA): similar to SPL

Experiment Results

- AV-WaN strongly outperforms all baselines and existing methods
- Our model improves the agent's efficiency substantially

	Heard sound				Unheard sounds					
Model	$\operatorname{SPL}\uparrow$	$\mathbf{SR}\uparrow$	NDG \downarrow	$NA\downarrow$	SNA↑	SPL ↑	$\mathbf{SR}\uparrow$	NDG \downarrow	$NA\downarrow$	SNA↑
Random Agent	4.7	18.5	89.7	448.9	1.7	4.7	18.5	89.7	448.9	1.7
Audio Direction Follower	54.7	72.0	19.5	173.0	41.1	11.1	17.2	55.4	415.2	8.4
Frontier Waypoints	44.0	63.9	47.5	218.5	35.2	6.5	14.8	96.3	434.2	5.1
Gan et al. [17]	57.6	83.1	17.6	130.5	47.9	7.5	15.7	76.3	336.6	5.7
Chen and Jain et al. [8]	75.0	95.7	1.3	84.4	43.4	28.6	41.2	16.0	223.1	14.9
AV-WaN (Ours)	85.3	98.0	0.6	38.8	70.0	54.4	71.2	12.3	173.6	43.8

Navigation Trajectories

- Gan et al.: is prone to errors and leads the agent to backtrack
- Chen et al.: oscillates around obstacles
- AV-WaN (ours): reaches the goal most efficiently



Idea: Semantic Audio Navigation

Sound informs navigating agent about...





Audio-Visual Navigation

Facebook AI Research Facebook Reality Labs University of Texas at Austin









Changan Unnat Chen Jain

it Se

Carl Sebastià V. Schissler Amengual Garí Sagnik Ziad Majumder Al-Halah











Santhosh Ruohan Vamsi Krishna Philip Ramakrishnan Gao Ithapu Robinson Kristen Grauman