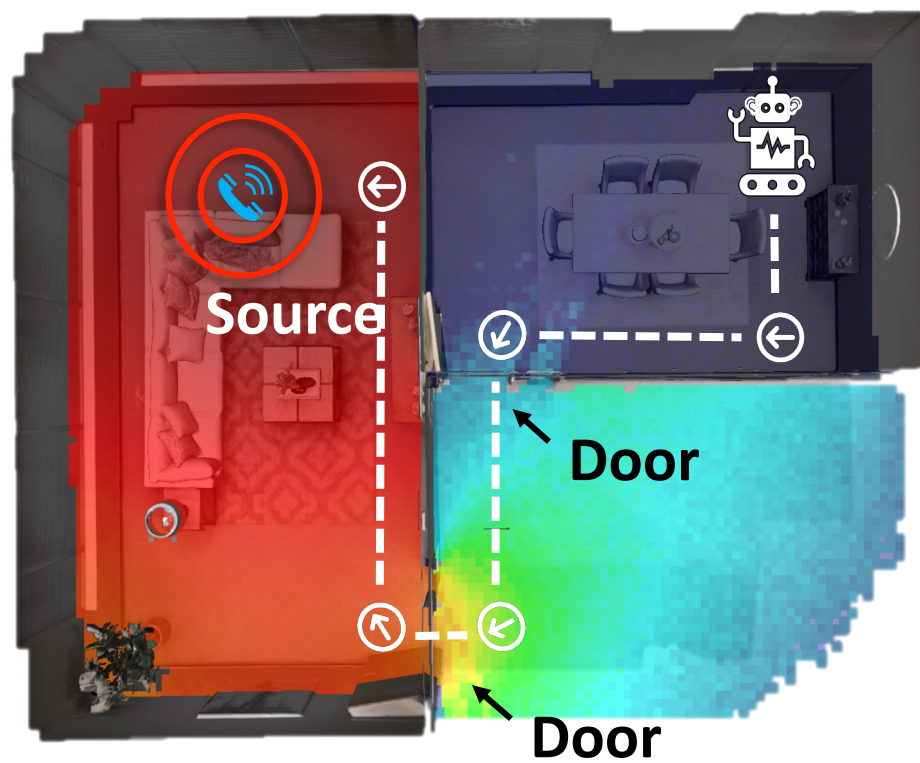


# SoundSpaces: Audio-Visual Navigation in 3D Environments

*Changan Chen<sup>\*1,4</sup>, Unnat Jain<sup>\*2,4</sup>, Carl Schissler<sup>3</sup>, Sebastia V. Amengual Gari<sup>3</sup>,  
Ziad Al-Halah<sup>1</sup>, Vamsi K. Ithapu<sup>3</sup>, Philip Robinson<sup>3</sup>, Kristen Grauman<sup>1,4</sup>*

*<sup>1</sup>UT Austin, <sup>2</sup>UIUC, <sup>3</sup>Facebook Reality Labs, <sup>4</sup>Facebook AI Research*



# Navigation Is a Multisensory Experience

We often use *vision*, *audio*, *touch*, *smell* to move around in the environment

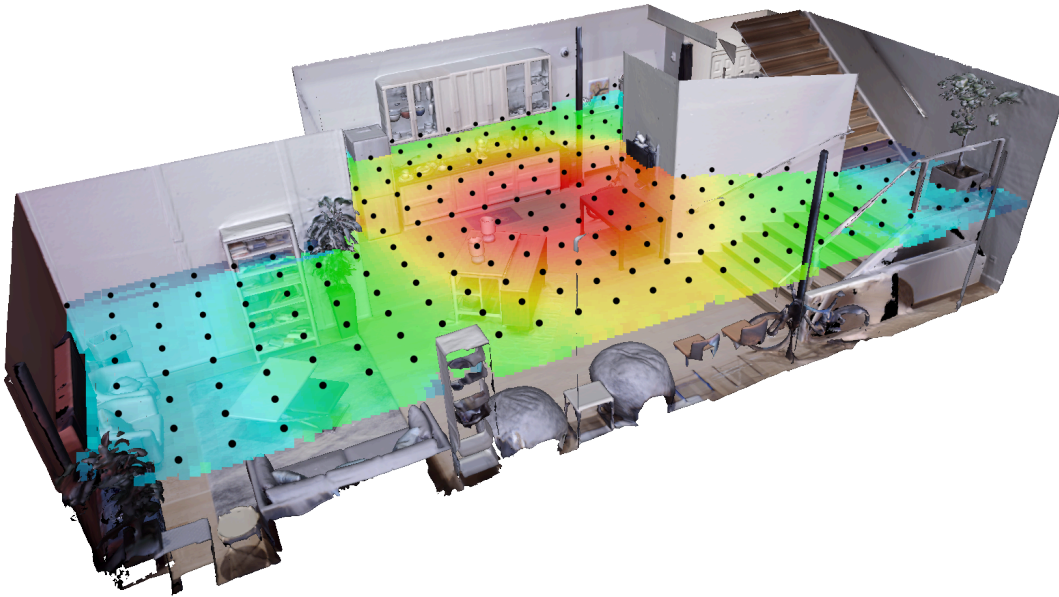
But today's embodied agent is **deaf**

We are the first to introduce audio-visual embodied navigation



# SoundSpaces: Our Audio Simulator

- We introduce SoundSpaces, an audio simulation platform to enable audio-visual navigation for two visually realistic 3D environments: Replica<sup>1</sup> and Matterport3D<sup>2</sup>



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

Table: Summary of dataset statistics

[1] The Replica Dataset: A Digital Replica of Indoor Spaces, Straub et al., arXiv, 2019

[2] Matterport3D: 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: Replica<sup>1</sup> and Matterport3D<sup>2</sup>
- 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

[1] The Replica Dataset: A Digital Replica of Indoor Spaces, Straub et al., arXiv, 2019

[2] Matterport3D: Learning from RGB-D Data in Indoor Environments, Chang et al., 3DV, 2017

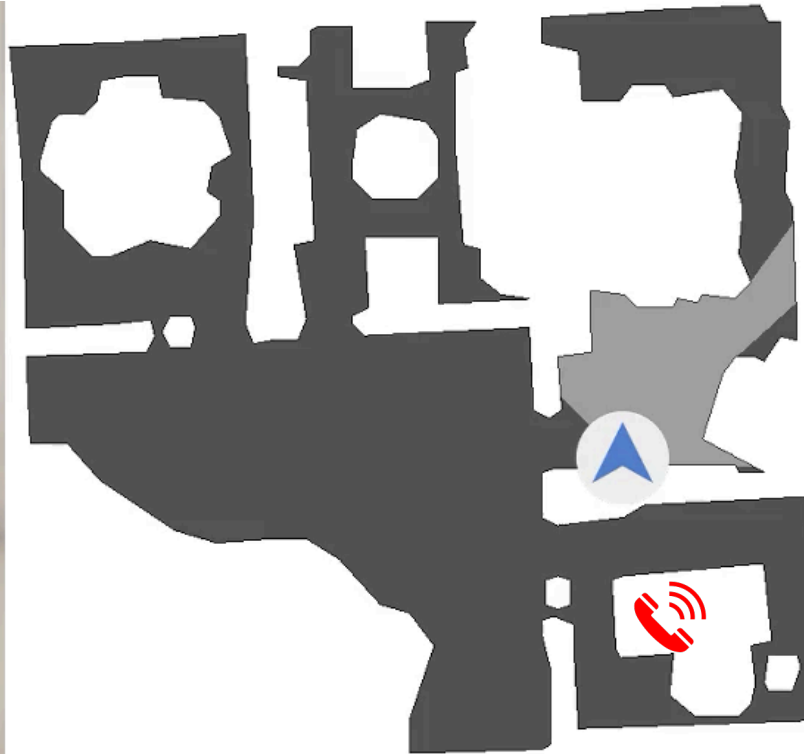


# Where Is My Phone?

Agent view



Top-down map

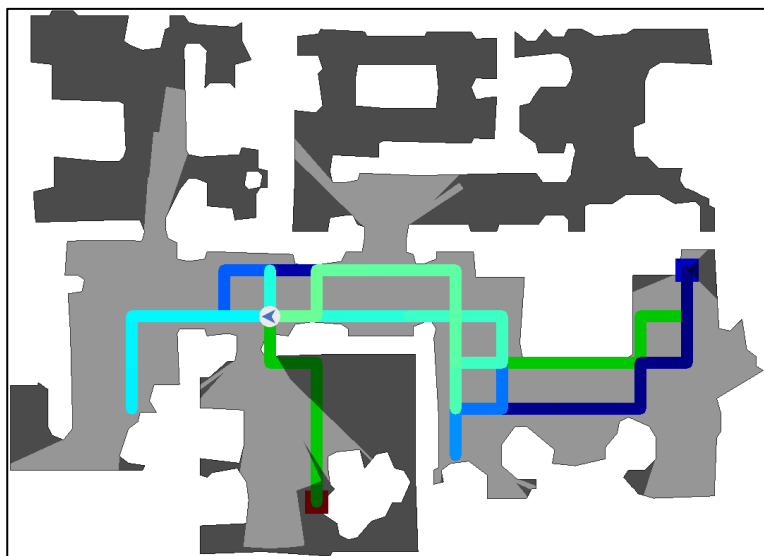


Direction: initially left ear is louder indicating sound coming from the left  
Intensity: overall intensity gets higher as the agent gets closer to the goal

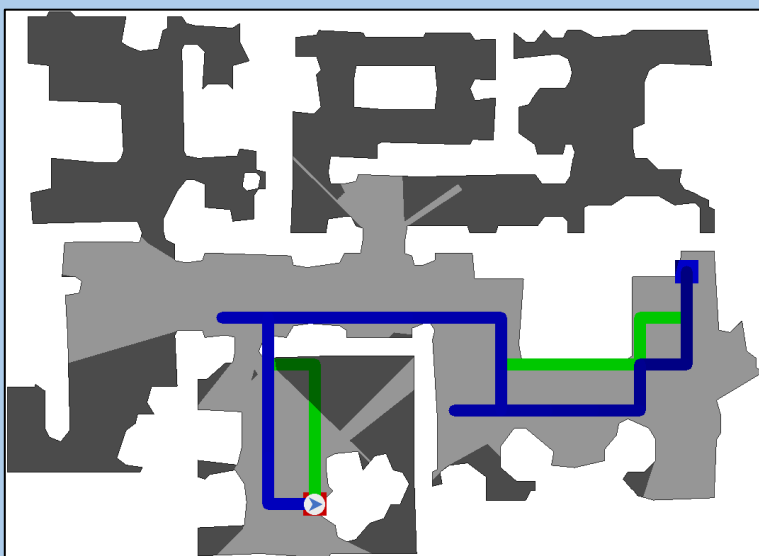
Agent Goal Start Shortest path Agent path Seen/Unseen area Occupied area

# Audio-Visual Navigation Tasks

PointGoal



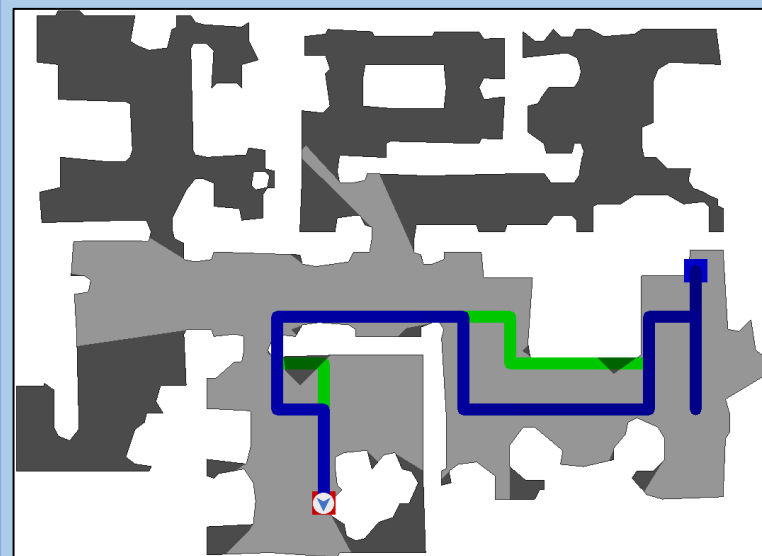
AudioGoal



AudioPointGoal



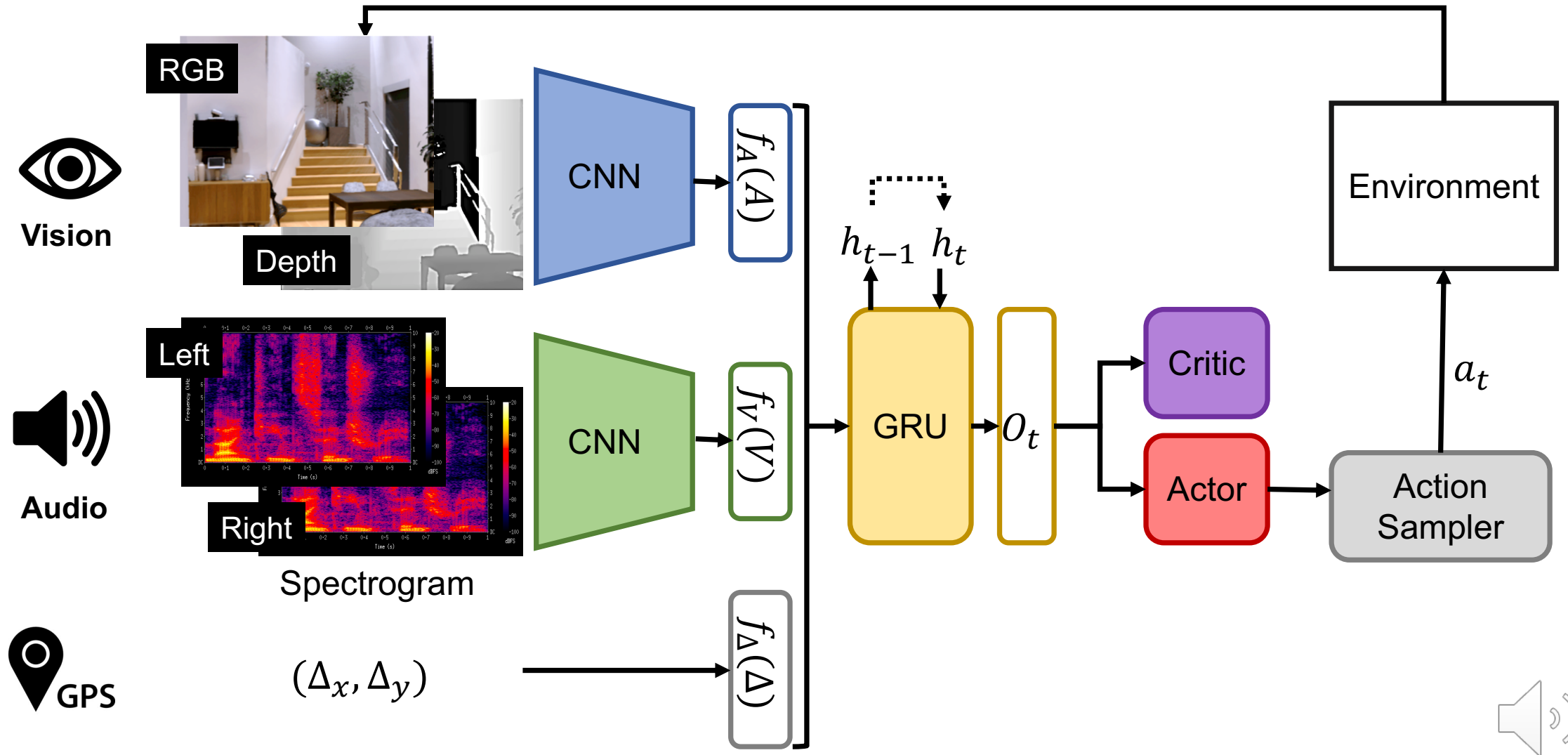
+



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



# Deep RL for Audio-Visual Navigation





# Experiment Results

We show:

- Audio helps navigation
- Audio supplants GPS for and audio target
- Our agent generalizes to unheard sounds

Table 3: Navigation performance (SPL) when generalizing to unheard sounds. Higher is better. Results are averaged over 7 test runs; all standard deviations are  $\leq 0.01$ .

Dataset		<i>PG</i>	<i>Same sound</i>		<i>Varied heard sounds</i>		<i>Varied unheard sounds</i>	
			<i>AG</i>	<i>APG</i>	<i>AG</i>	<i>APG</i>	<i>AG</i>	<i>APG</i>
Replica	Blind	0.480	0.673	0.681	0.449	0.633	0.277	0.649
	RGB	0.521	0.626	0.632	0.624	0.606	0.339	0.562
	Depth	0.601	0.756	0.709	0.645	0.724	0.454	0.707
Matterport3D	Blind	0.426	0.438	0.473	0.352	0.500	0.278	0.497
	RGB	0.466	0.479	0.521	0.422	0.480	0.314	0.448
	Depth	0.541	0.552	0.581	0.448	0.570	0.338	0.538





# SoundSpaces: Audio-Visual Navigation in 3D Environments

*Changan Chen<sup>\*1,4</sup>, Unnat Jain<sup>\*2,4</sup>, Carl Schissler<sup>3</sup>, Sebastia V. Amengual Gari<sup>3</sup>,  
Ziad Al-Halah<sup>1</sup>, Vamsi K. Ithapu<sup>3</sup>, Philip Robinson<sup>3</sup>, Kristen Grauman<sup>1,4</sup>*

*<sup>1</sup>UT Austin, <sup>2</sup>UIUC, <sup>3</sup>Facebook Reality Labs, <sup>4</sup>Facebook AI Research*

For more details, please refer to our paper:

<https://arxiv.org/pdf/1912.11474.pdf>

For more simulation demo, please check our project page:

[http://vision.cs.utexas.edu/projects/audio\\_visual\\_navigation](http://vision.cs.utexas.edu/projects/audio_visual_navigation)

