# Semantic Audio-Visual Navigation

Changan Chen<sup>1,2</sup>, Ziad Al-Halah<sup>1</sup>, Kristen Grauman<sup>1,2</sup> <sup>1</sup>UT Austin,<sup>2</sup>Facebook Al Research

**CVPR 2021** 







### **Audio-Visual Navigation**

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

- Travel to an unknown sounding object in an unmapped environment
- Receive both visual and auditory observations





#### 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



# Semantic AudioGoal Task



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 SoundSpaces<sup>1</sup> with semantic sounds
- 21 object categories in Matterport3D<sup>2</sup>: chair, TV, cabinet, sink etc.
- Object-emitted sounds and object-related sounds



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

### Our Idea

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



# Semantic Audio-Visual Navigation (SAVi)



# **Navigation Results**

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



# Navigation Example



The sound stops before the agent reaches the goal, but the agent identifies its chest-of-drawers sound and uses its visual perception to locate the target object/

# Semantic Audio-Visual Navigation

Changan Chen<sup>1,2</sup>, Ziad Al-Halah<sup>1</sup>, Kristen Grauman<sup>1,2</sup> <sup>1</sup>UT Austin,<sup>2</sup>Facebook Al Research

CVPR 2021

Code and audio simulation data available at: <u>http://vision.cs.utexas.edu/projects/semantic-audio-visual-navigation</u>

SoundSpaces Challenge at Embodied AI Workshop, CVPR 2021: <u>https://soundspaces.org/challenge</u>



