# **Compositional Audio Representation Learning**

# Sripathi Sridhar, Mark Cartwright

New Jersey Institute of Technology

**Correspondence**: <u>ss645@njit.edu</u>

# Highlights

- We propose a source-centric representation learning framework for environmental sound
- We investigate the design space of these models
- Compositional audio representations outperform baselines on audio classification
- Label supervision helps learn better sourcecentric representations on synthetic data
- Proposed additional loss terms lead to better downstream performance

# **SINC** AND COMPUTING LAP

### Motivation

- flexibility and interpretability

### What is Compositional Audio Representation Learning?

**Goal**: Map each source in the auditory scene to a distinct embedding, taking inspiration from object-centric learning in computer vision.



Source-centric embedding



### Synthetic data

Open-set Soundscapes · 260k 10s (OSS)

Open-set Tagging (OST)

- soundscapes synthesized from FSD50K sources
- 54 seen classes, 35 unseen classes
- 500k 1s clips Window around center of each event in soundscape

### **Research questions**

- models?
- classes?

## **Key findings**

- **1.** Compositional audio representations outperform baselines and are useful for environmental sound classification
- 2. Target signal matters! Reconstructing features is significantly better than reconstructing spectrograms
- 3. The proposed disjointedness and sparsity loss terms lead to better downstream performance
- 4. Label supervision helps learn better source-centric representations on synthetic data

MLP feat decoder

CNN spec decoder



Tfmer w/o disjo	
MLP	
MLF	del
Tfmer w/	Mo
Tfmer w/o recor	

