

TriCycle: Audio Representation Learning from Sensor Network Data Using Self-Supervision

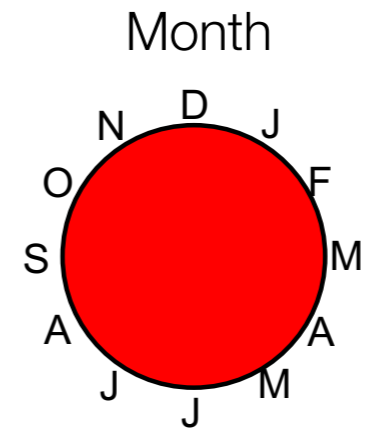
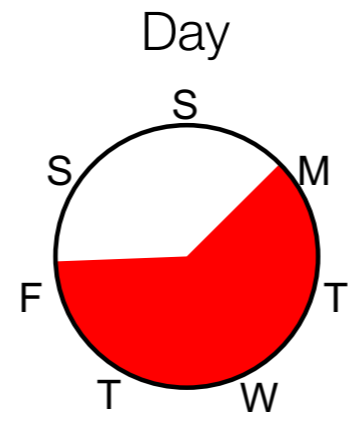
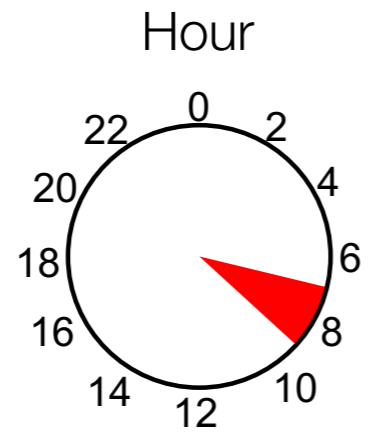
Mark Cartwright¹, Jason Cramer¹, Justin Salamon², Juan Pablo Bello¹

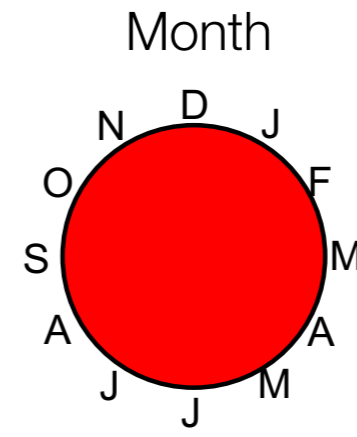
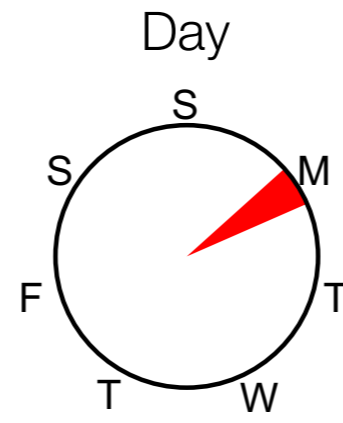
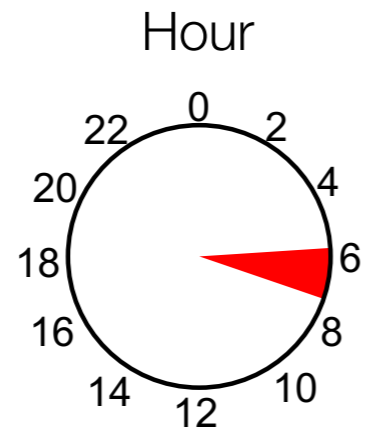
1. New York University Music and Audio Research Lab

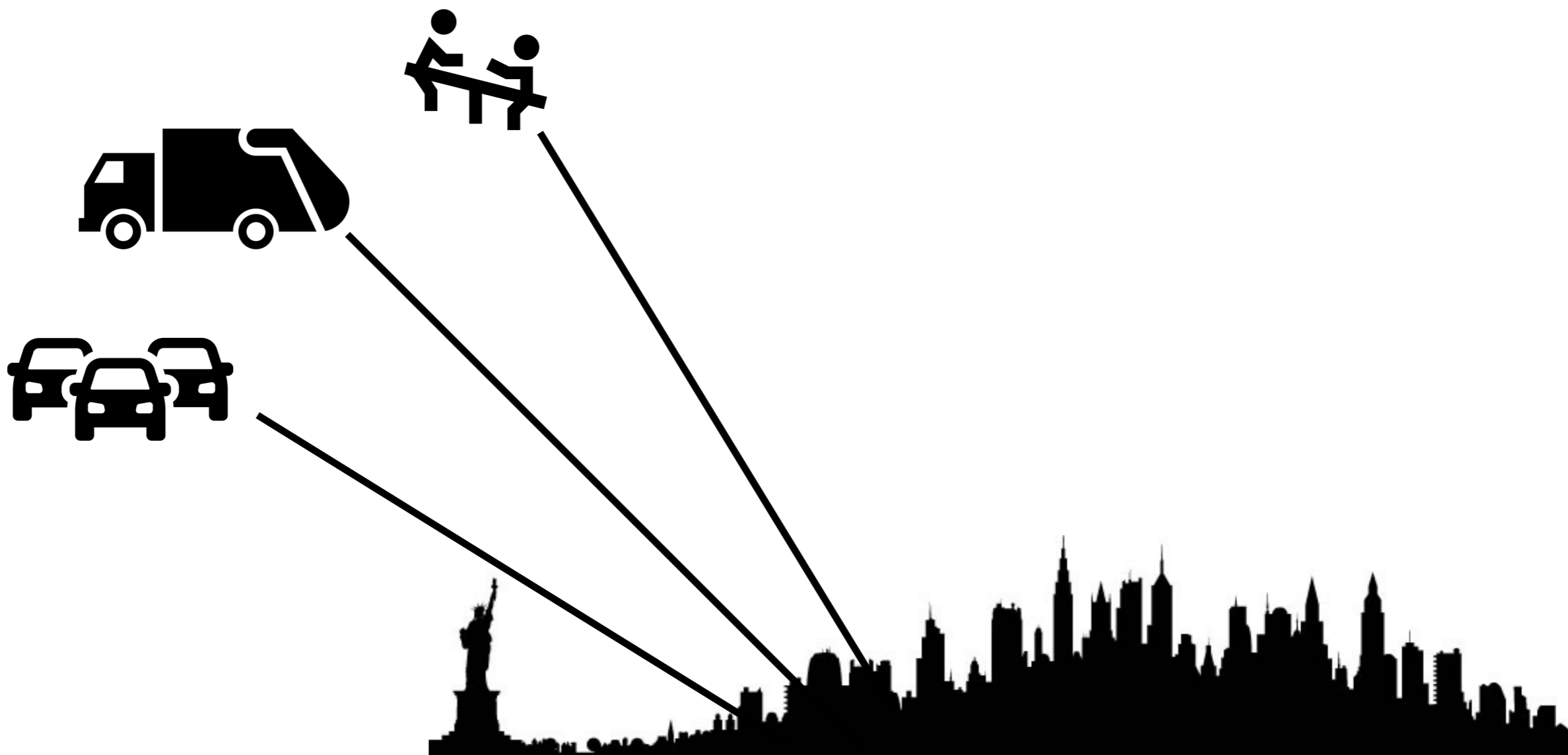
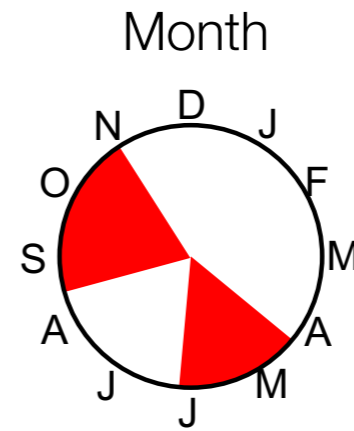
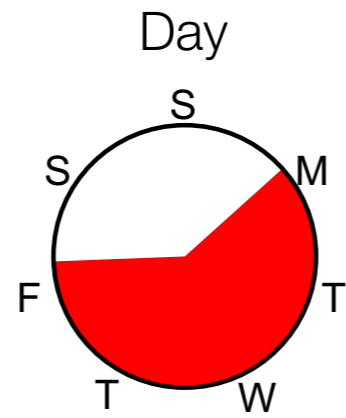
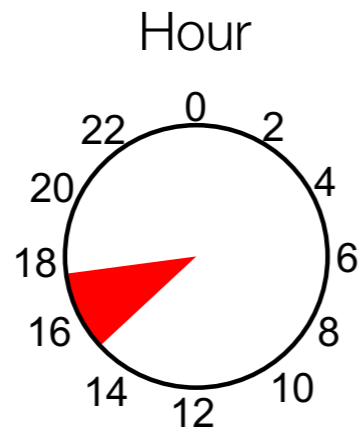
2. Adobe Research

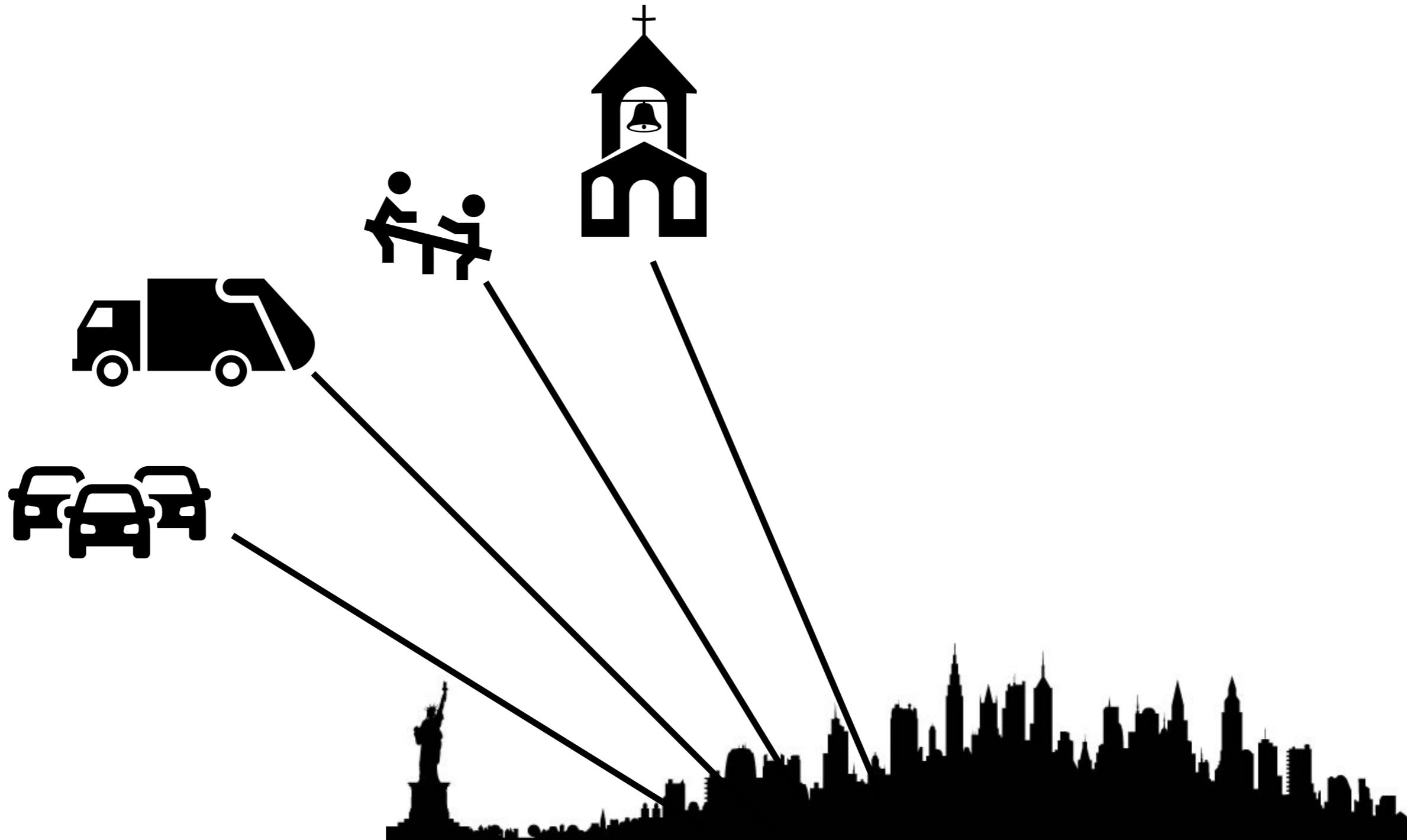
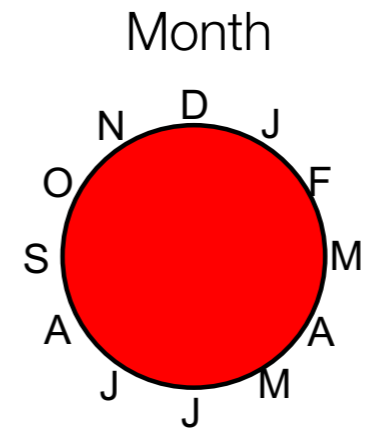
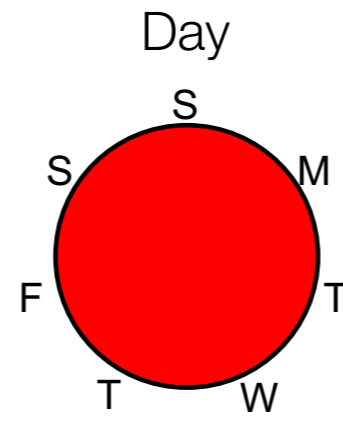
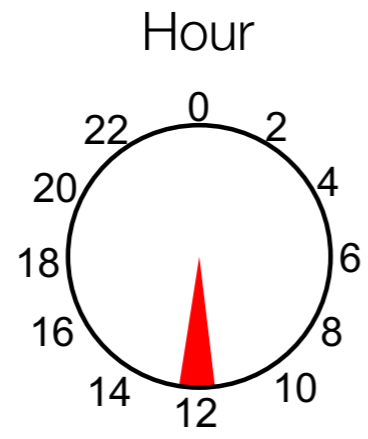


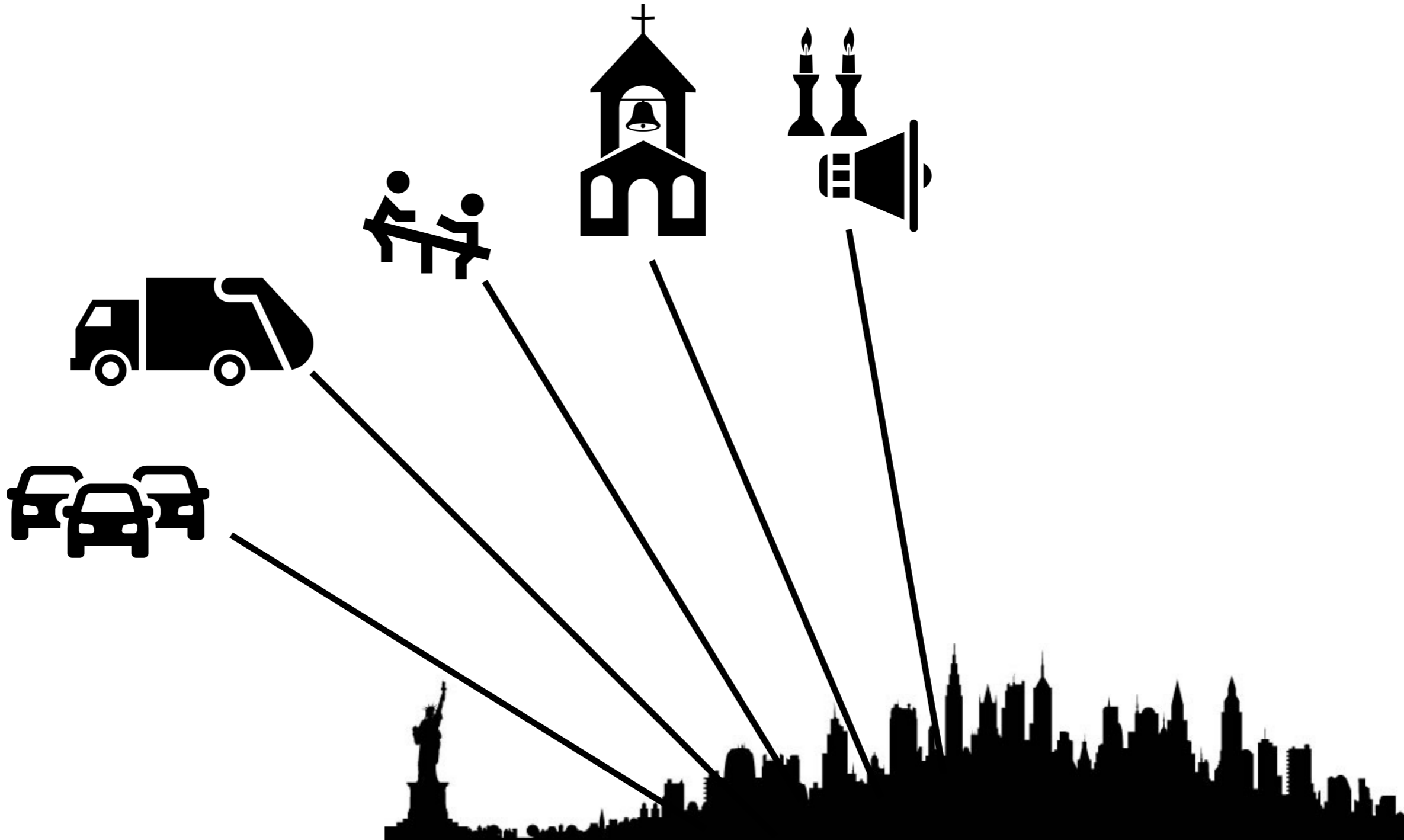
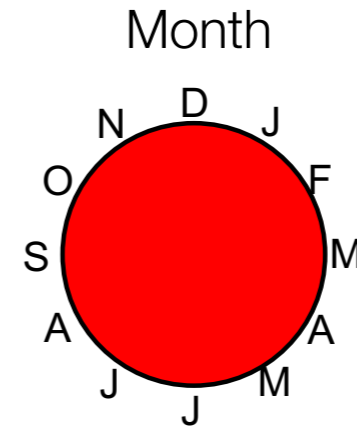
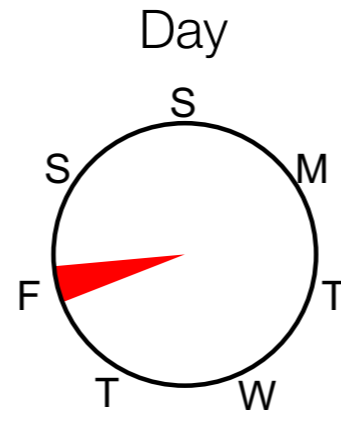
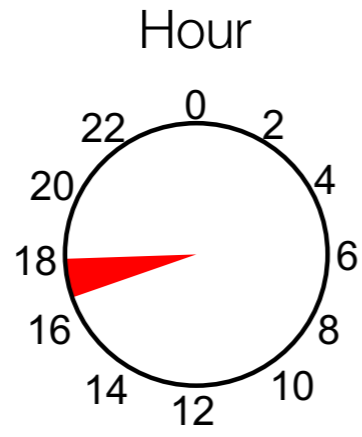


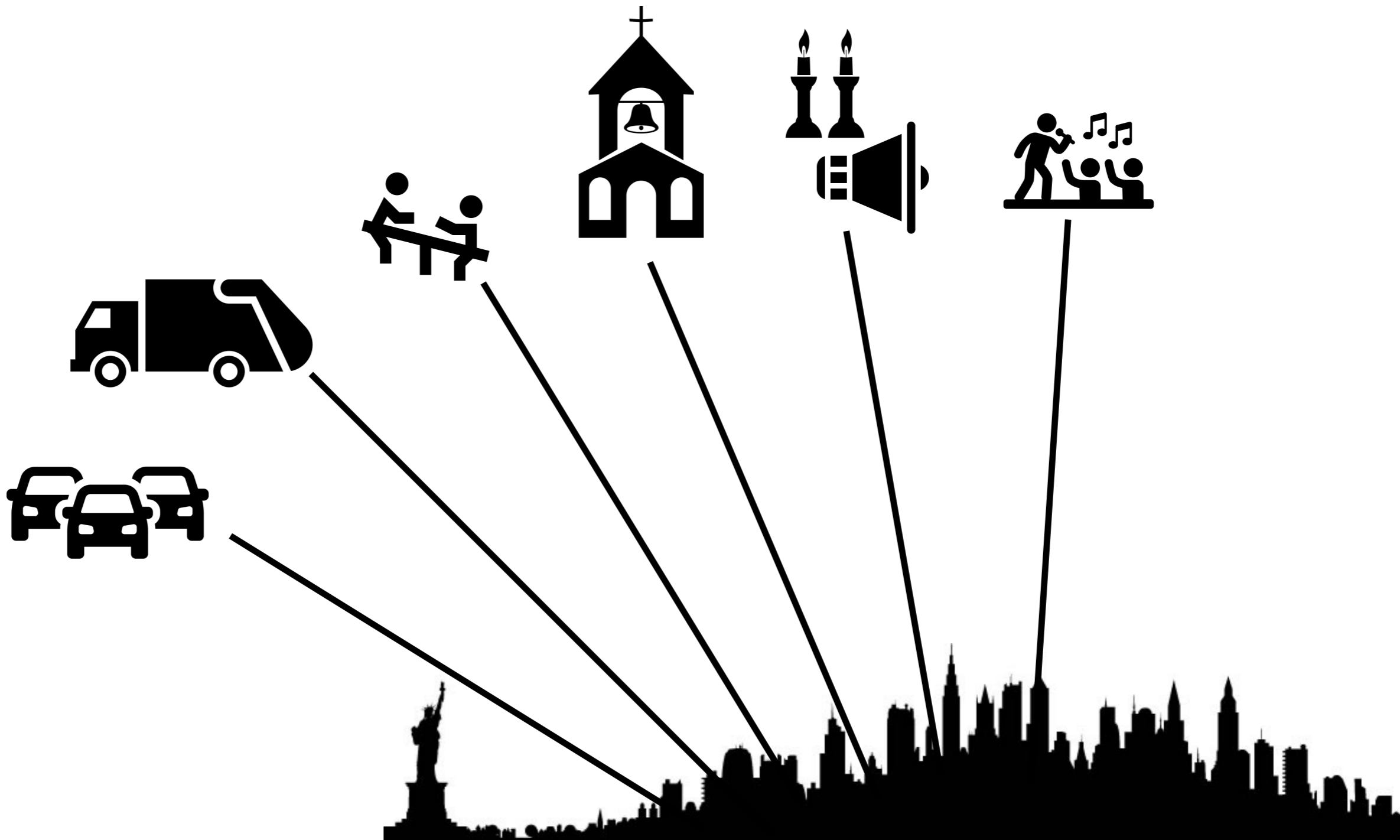
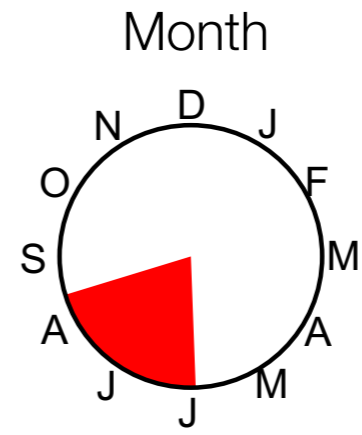
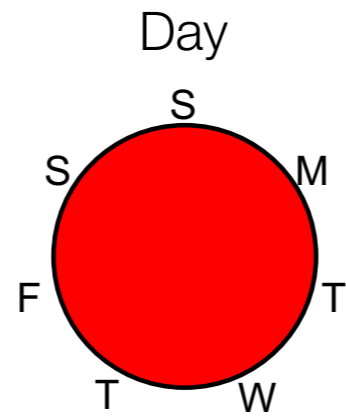
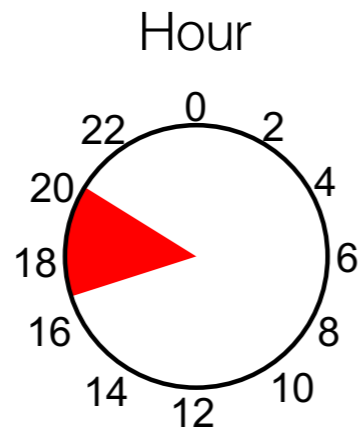


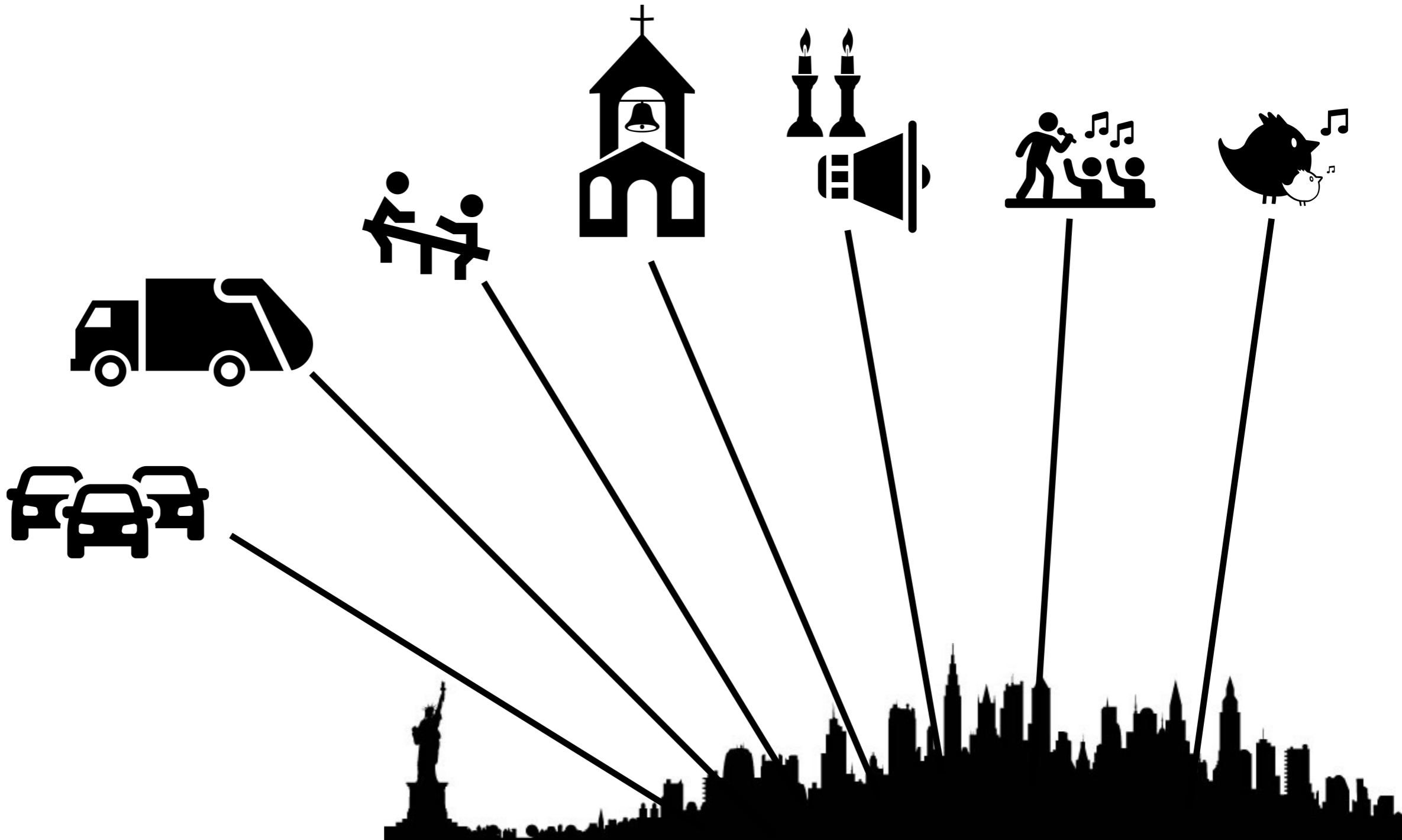
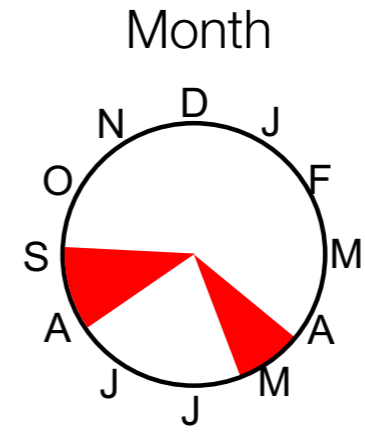
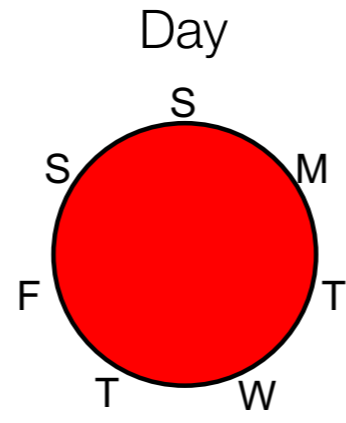
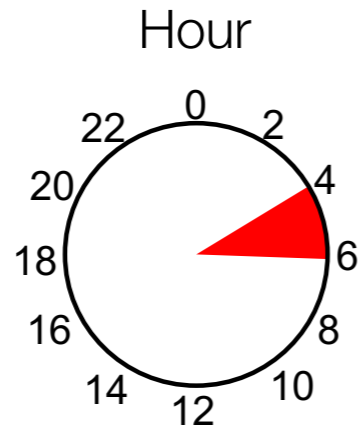














57

Fixed-location Sensors

130M

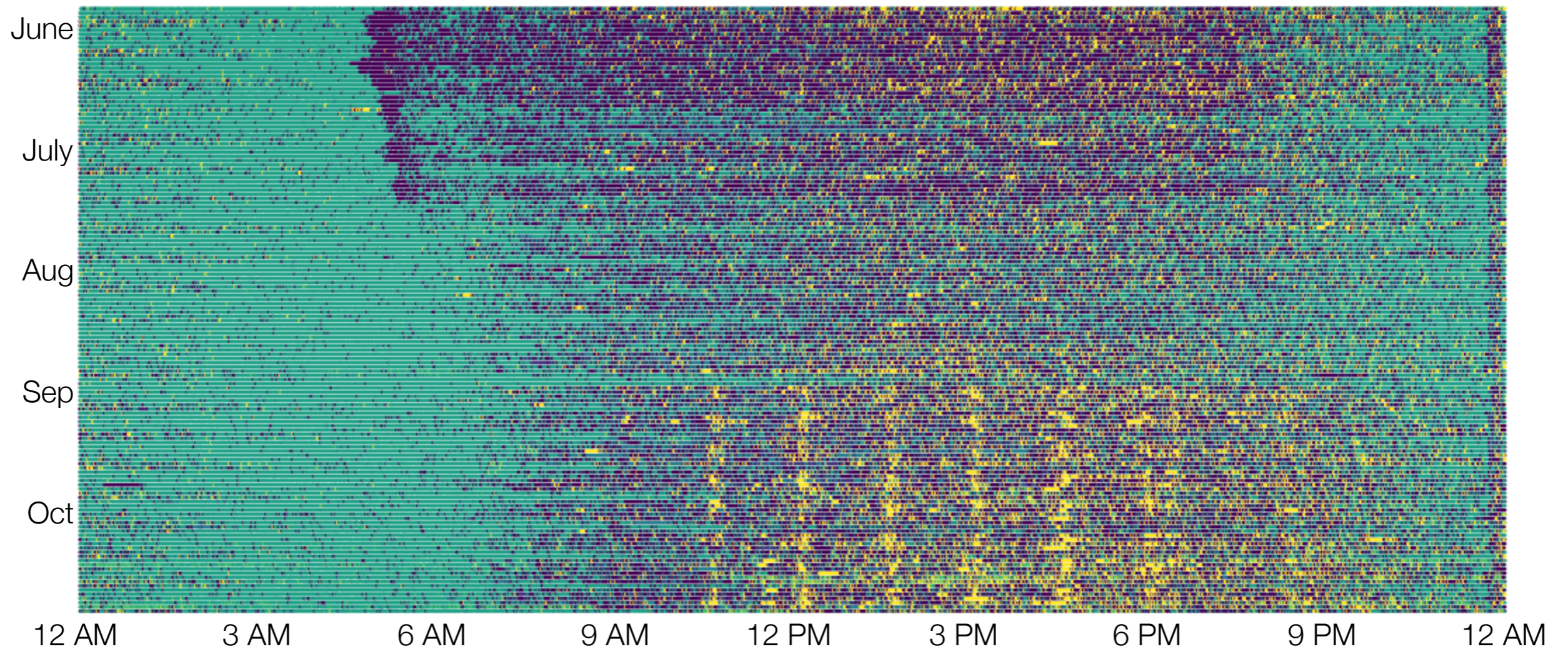
10 sec Recordings

40

Sensor Years of Audio

Long-term temporal structure in SONYC recordings

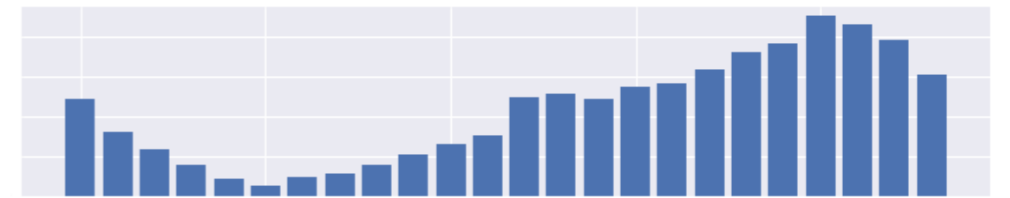
Predominant cluster (N=8) over 4 months for 1 sensor



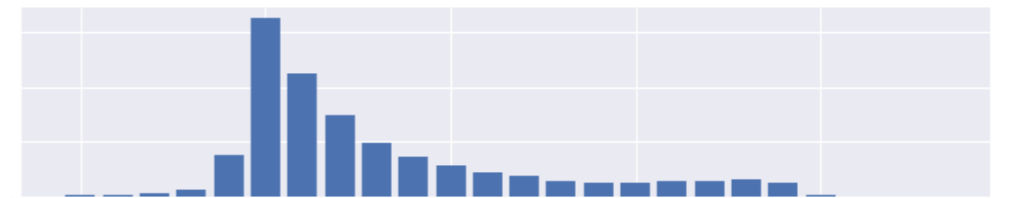
Long-term temporal structure in SONYC recordings

Cluster frequency grouped by hour of the day on SONYC recordings

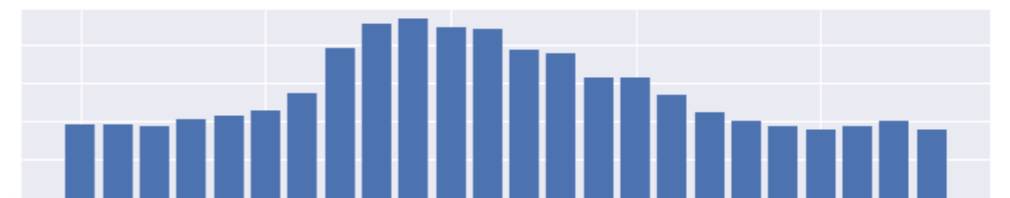
“People Talking”



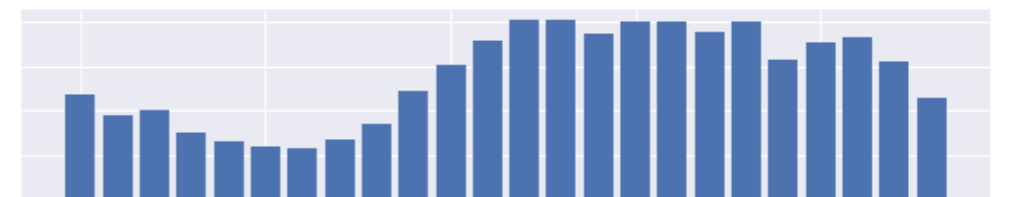
“Birds Chirping”



“Loud Truck Engine”



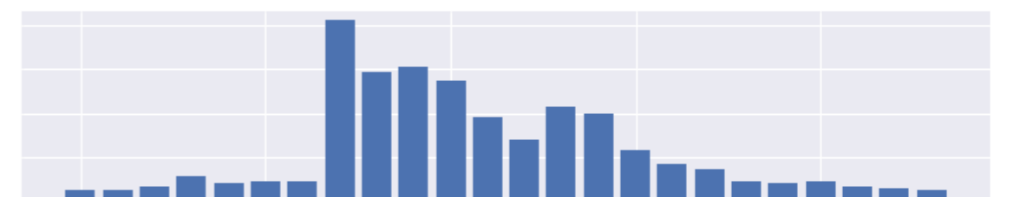
“Siren”



“Dog Barking”



“Powered Saw”

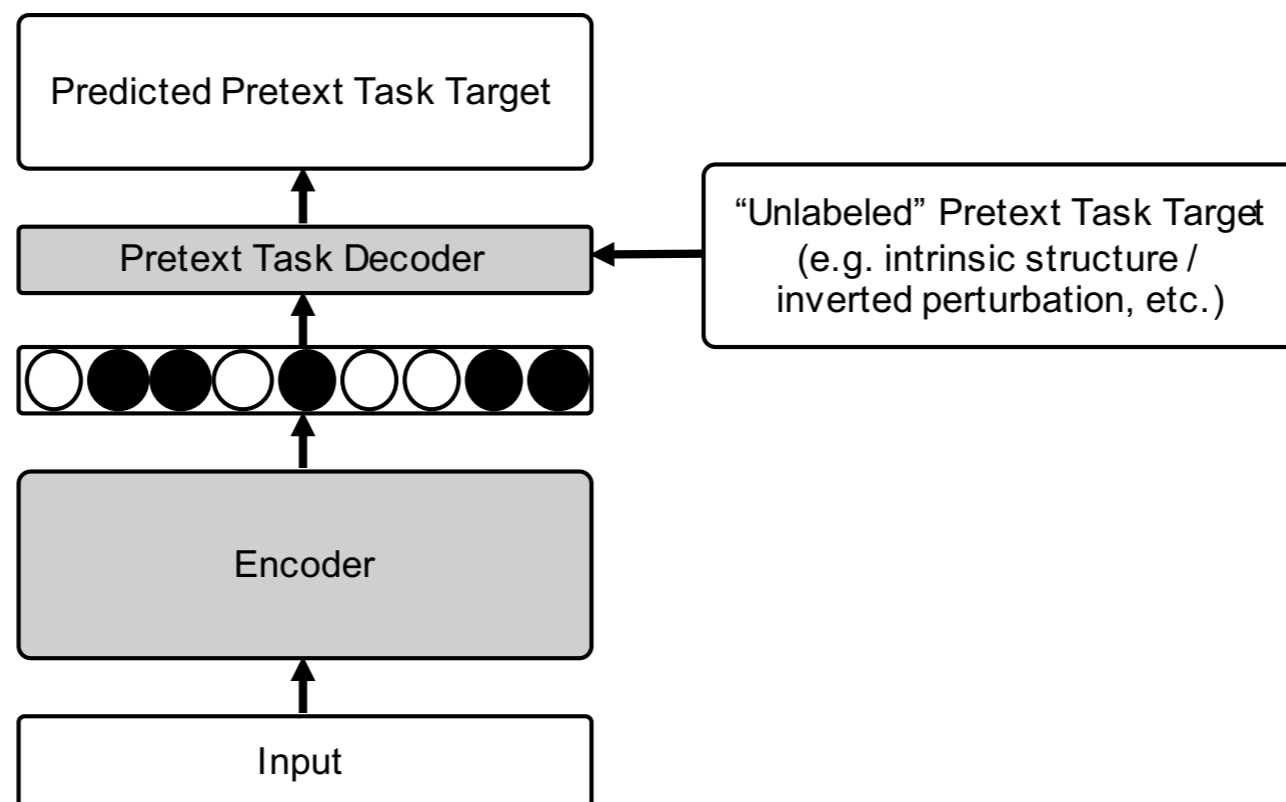


12 AM 6 AM 12 PM 6 PM 12 AM
Hour

Can we exploit this long-term seasonal structure for self-supervised audio representation learning?

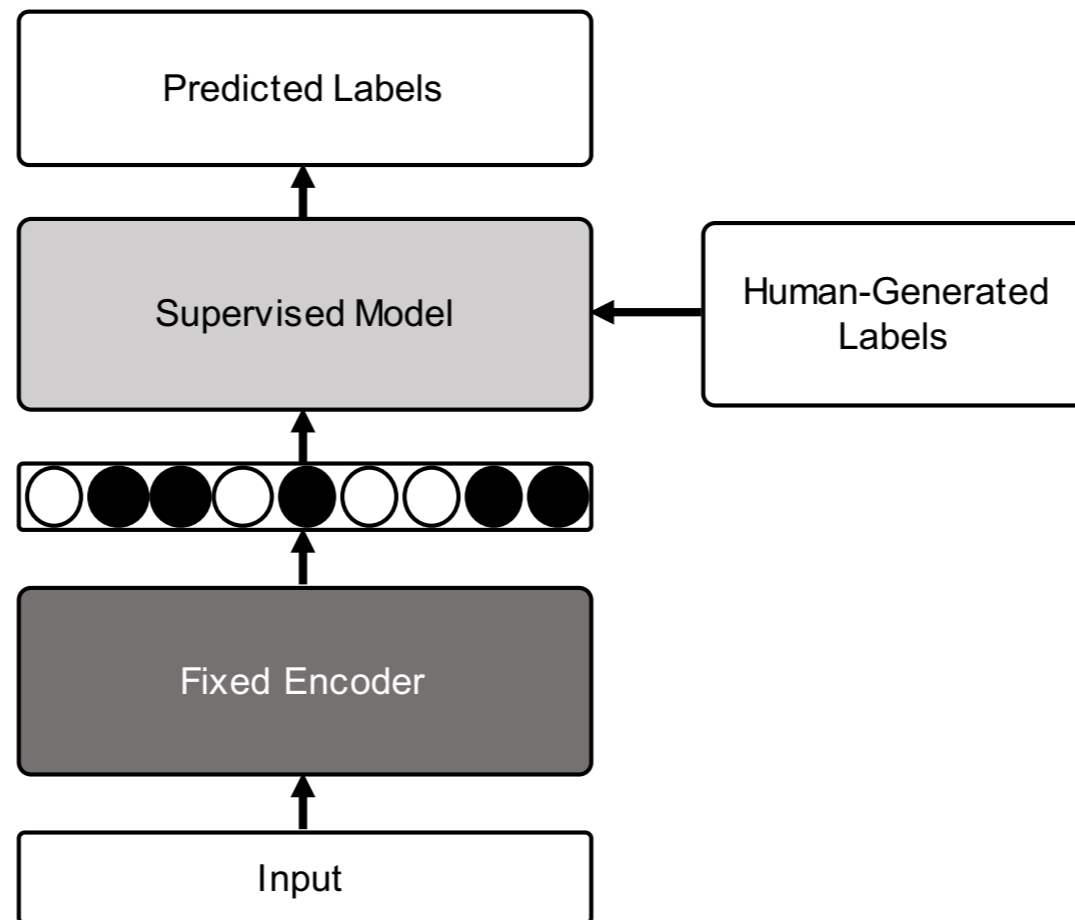
Self-supervised pretext task

- Learn representations (embeddings) by solving pretext tasks
- Pretext tasks exploit known intrinsic structure or estimate / invert a controlled perturbation
- Key is that pretext tasks **do not** require (human-generated) labels and are trained on **lots** of this “unlabeled” data



Supervised downstream task

- With learned representation as input, use **simpler, smaller** capacity supervised model with **fewer labeled examples** in a downstream task



Examples in computer vision

Unsupervised Visual Representation Learning by Context Prediction

¹ Scho
Carné

Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles

UNSUPERVISED REPRESENTATION LEARNING BY PREDICTING IMAGE ROTATIONS

Spyros
Univers
Ecole d
{spyro



Image

Learning and Using the Arrow of Time

Donglai Wei¹, Joseph Lim², Andrew Zisserman³ and William T. Freeman^{4,5}

¹Harvard University ²University of Southern California

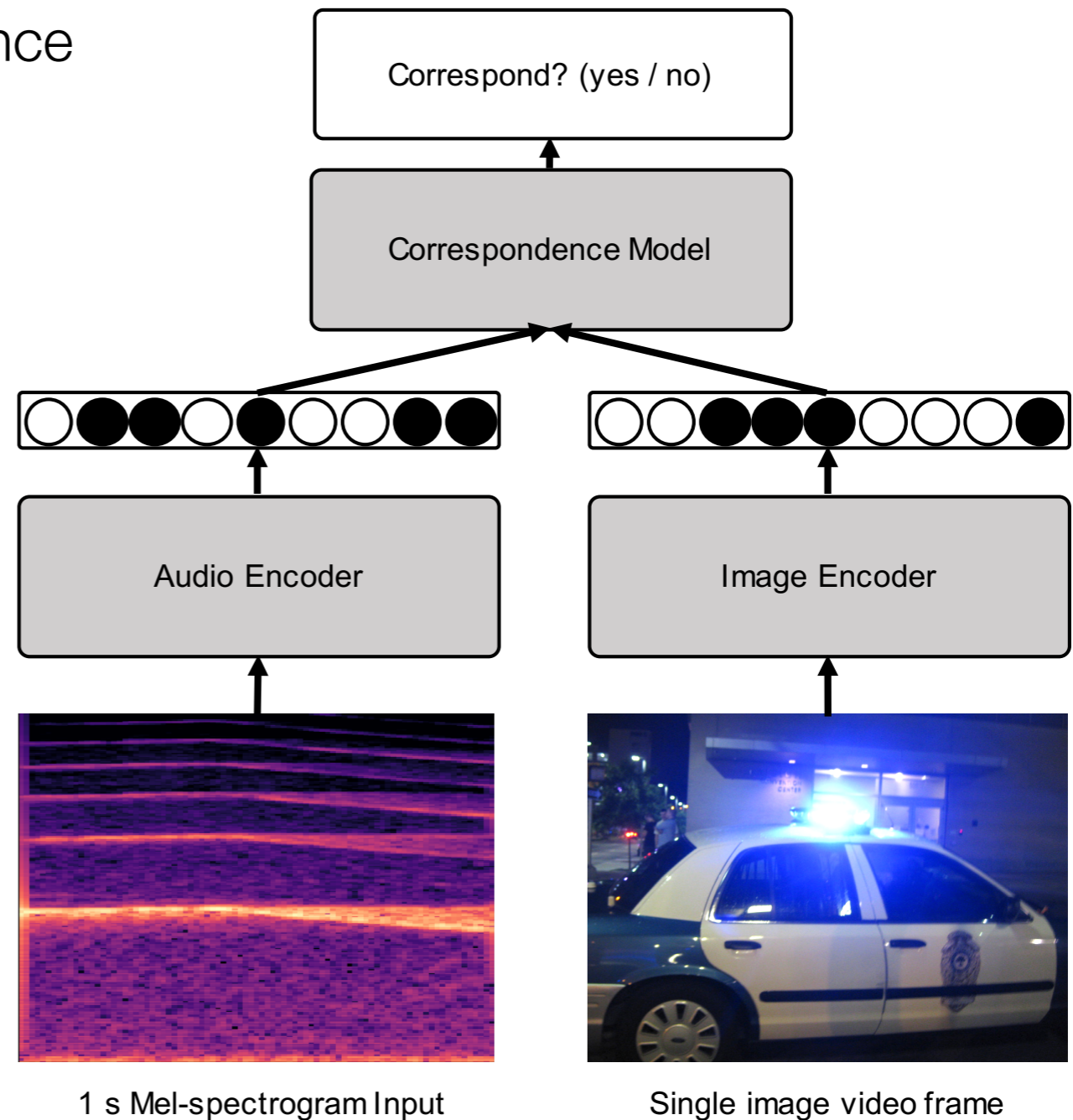
³University of Oxford ⁴Massachusetts Institute of Technology ⁵Google Research

donglai@seas.harvard.edu, limjj@usc.edu, az@robots.ox.ac.uk, billf@mit.edu



Examples in machine listening

- Arandjelovic & Zisserman, “Look, listen and learn” (L3), ICCV 2017
- Exploits audio-visual correspondence

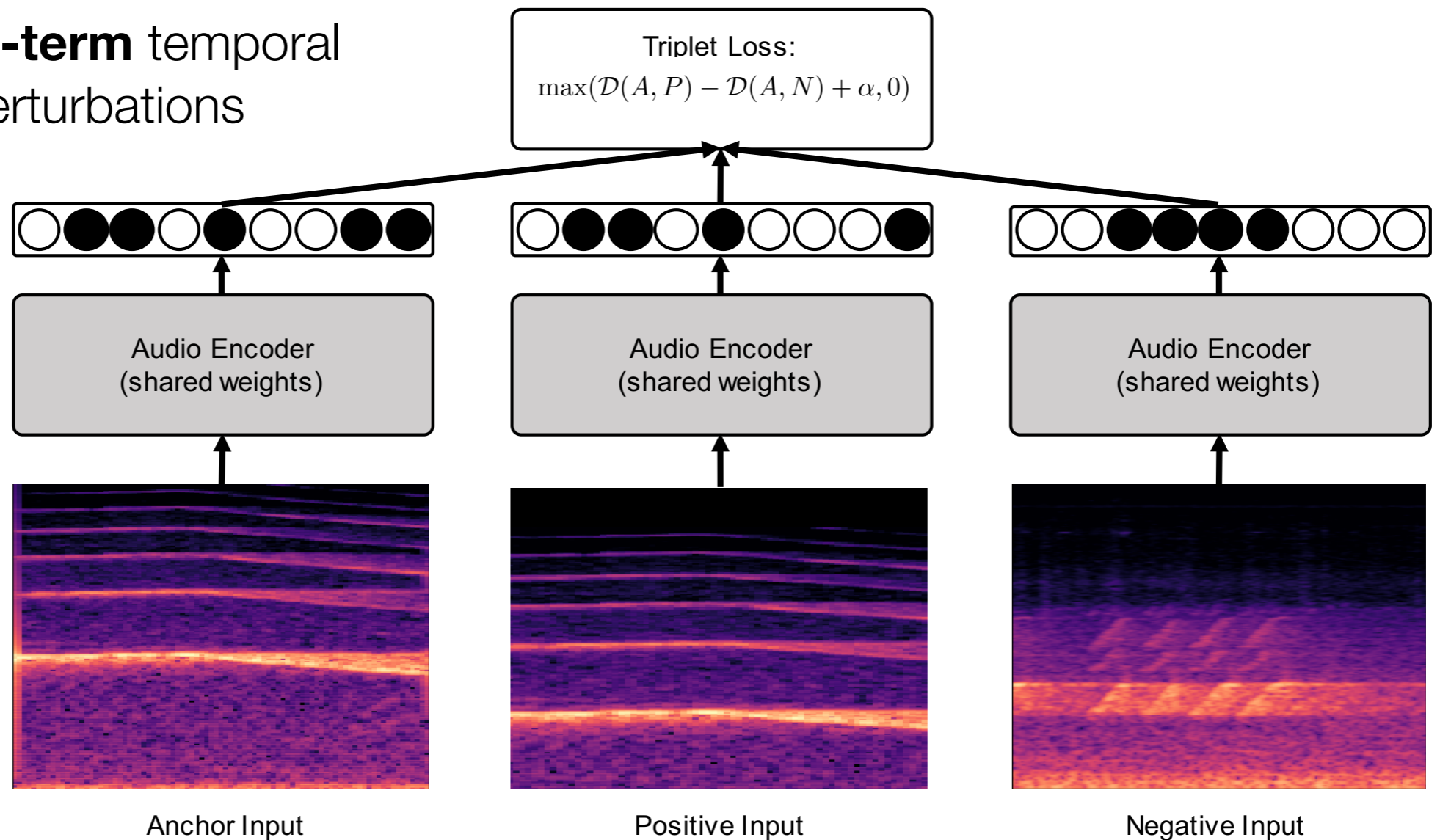


1 s Mel-spectrogram Input

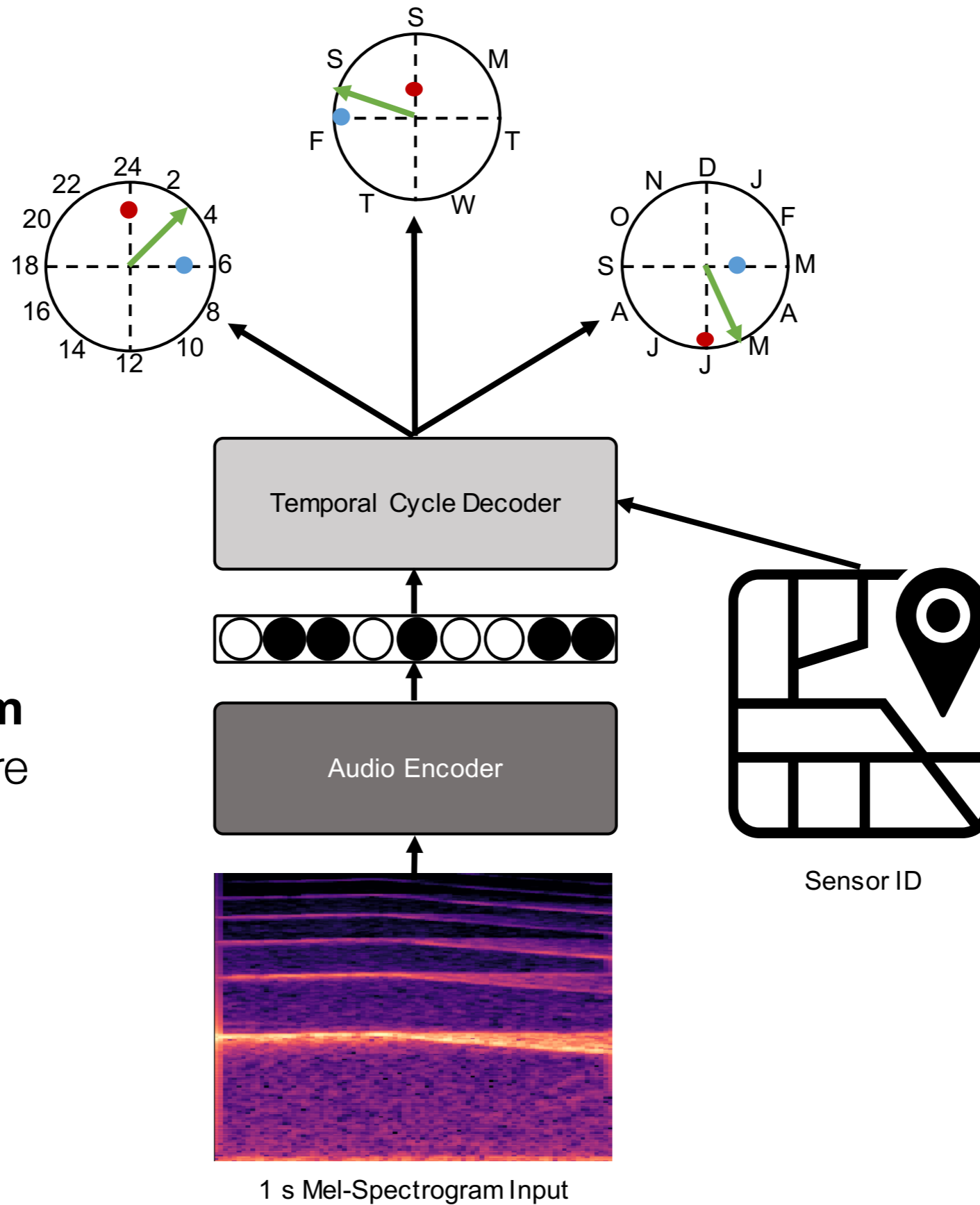
Single image video frame

Examples in machine listening

- Jansen, et al. “Unsupervised learning of semantic audio representations”, ICASSP 2018
- Exploits **short-term** temporal structure or perturbations



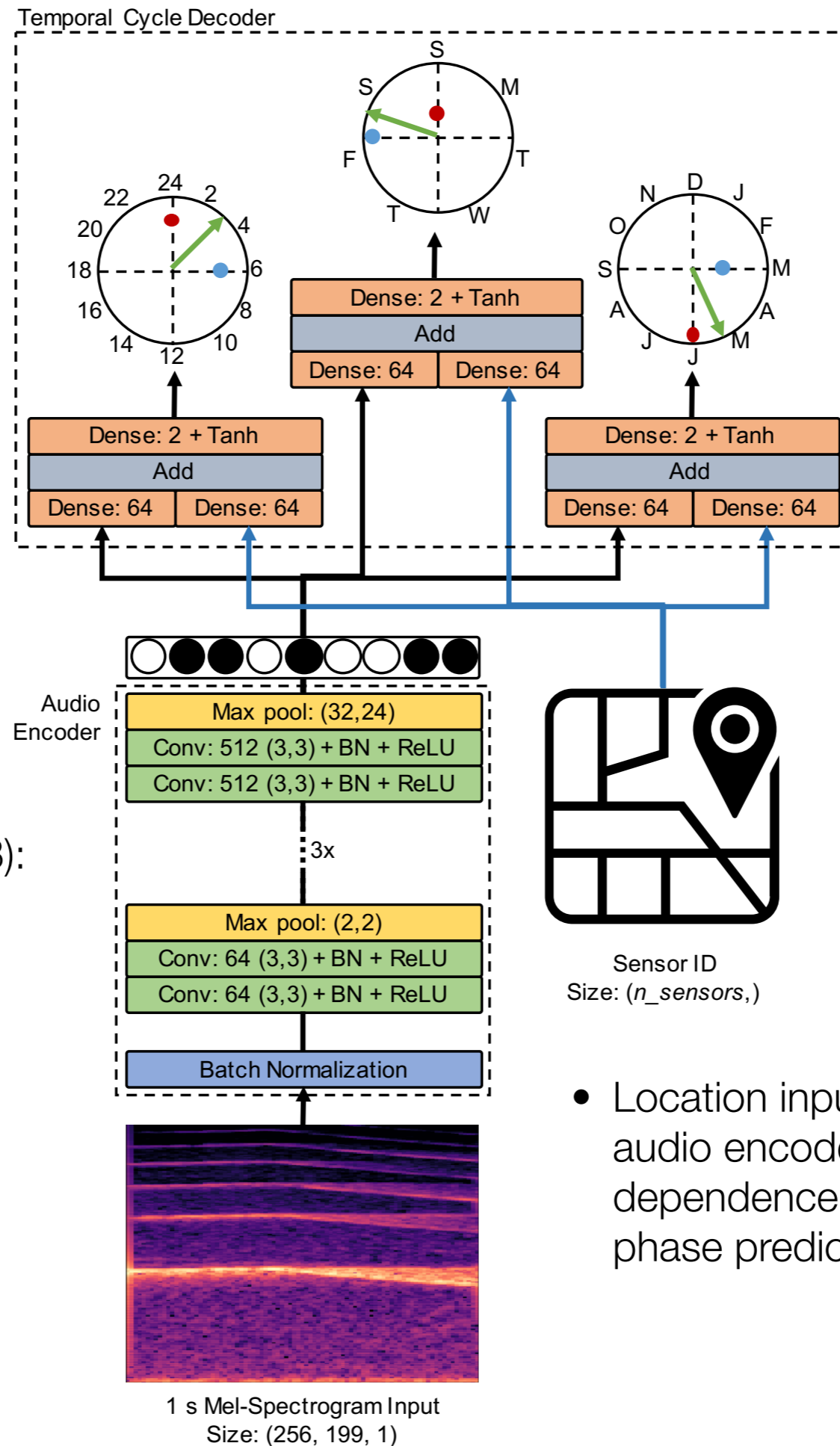
TriCycle Model



- We propose to exploit **long-term** temporal structure

TriCycle Model

- Audio encoder same as “Look, Listen, and Learn” (L3):
 - Simple CNN
 - 4 convolutional blocks
 - Each with 2 conv. layers + max pooling
- Input:
 - 48kHz
 - 256-bin Mel spectrogram
 - log-scaled magnitude
 - 5 ms hop



- To avoid issues with phrase wrapping, phase encoded as $[\cos(\phi), \sin(\phi)]$ optimized with MSE loss

- Location input incorporated after audio encoder to account for location dependence of sound events in phase prediction

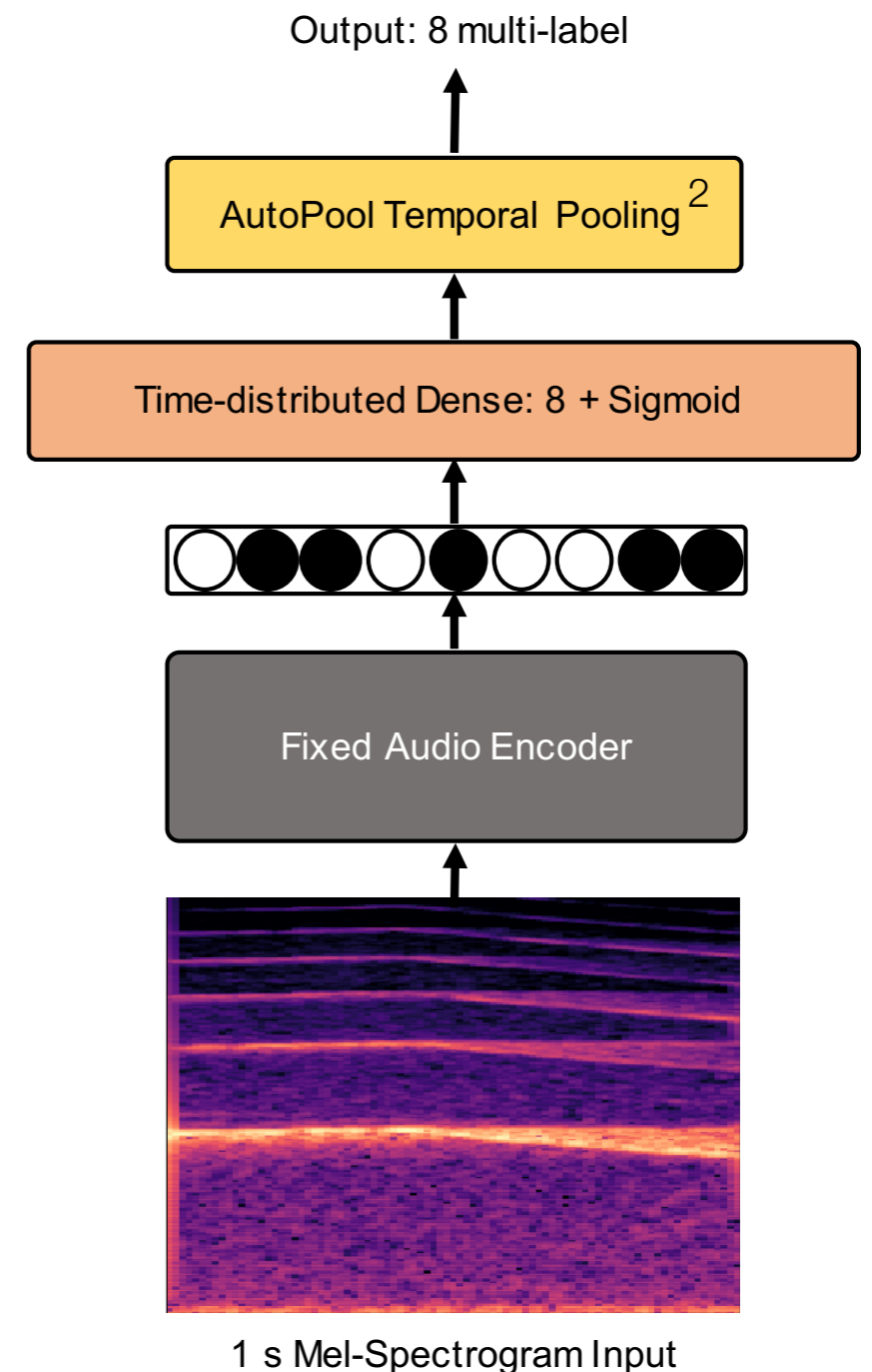
TriCycle Training

- Because of resource constraints, limited SONYC dataset to 2017 data from 25 sensors ~25M 10 sec recordings (69k hours)
- Randomly sampled
- 1500 “epochs” (24M training examples)

Supervised downstream task: Urban sound tagging

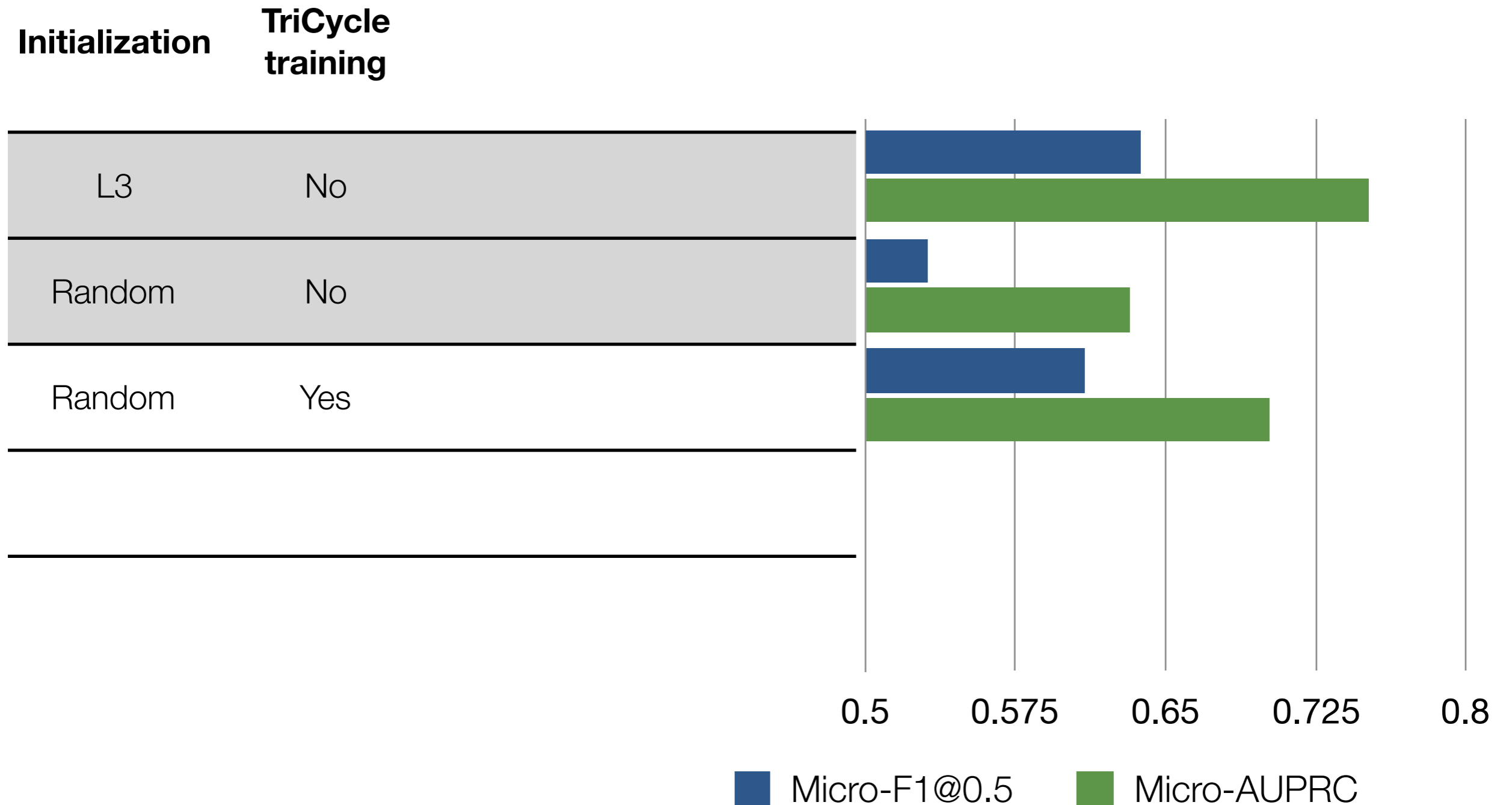
SONYC Urban Sound Tagging (UST) Dataset¹

- labeled subset of SONYC data
- v0.1 Released in March
- 2019 DCASE Urban Sound Tagging Challenge dataset
- 10 sec recordings from SONYC sensors
2351 training
443 validation
274 test (*did not use*)
- Weak multi-label annotation on 23 fine-level classes from 8 coarse-level groups (we used the coarse labels):
engine, machinery impact, non-machinery impact, powered saw, alert signal, music, human voice, dog
- 3 Zooniverse volunteer annotators per recording
Used minority vote to aggregate
- Validation and test set annotated by SONYC team



1. Cartwright, et al. "SONYC Urban Sound Tagging (SONYC-UST): A multilabel dataset from an urban acoustic sensor network", DCASE 2019
2. McFee, Salamon, Bello. "Adaptive pooling operators for weakly labeled sound event detection", TASLP 2018

Urban sound tagging results with TriCycle



Strategies to focus on foreground events: High-activity sampling

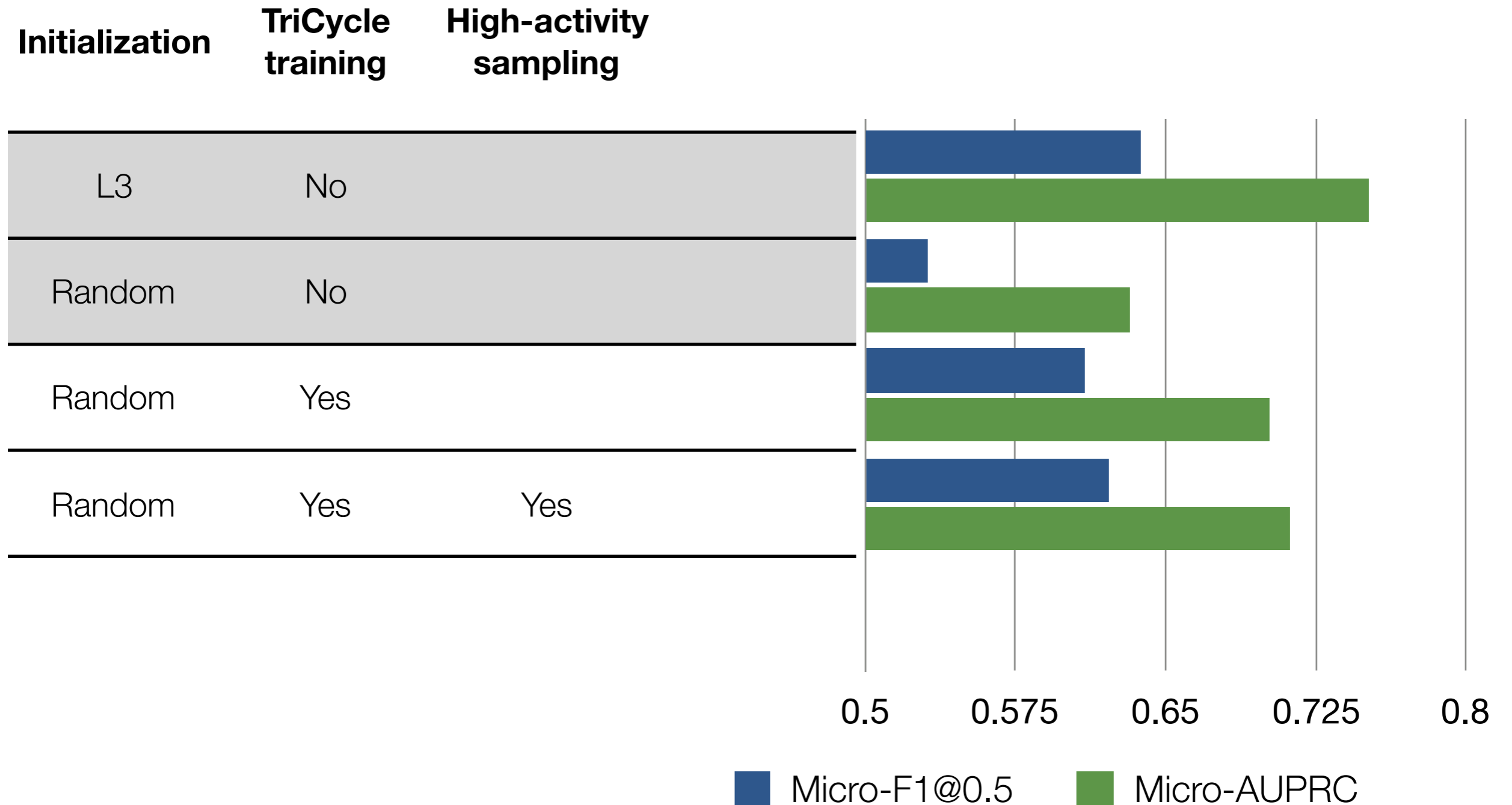
- Focus on high activity regions but still evenly sample each hour
- Compute SPL “activity” metric for each 10 s recording (SPL b/c precomputed):

$$\sqrt{\sum_{n=0}^{79} (d_{m,n} - d_{m,n-1})^2}$$

for SPL sequence d of length 80 (i.e., 10 s with 0.125 s step size) from sensor m

- Only sample from top 15 percent of each hour
- Within each 10 s recording, sample 1 s clip, weighting by SPL

Urban sound tagging results with TriCycle



Focusing on foreground events: Per-Channel Energy Normalization (PCEN)

Pre-process with Per-Channel Energy Normalization (PCEN)¹

- Spectrogram processing that Gaussianizes and decorrelates frequency channels while retaining sound events of interest (parameter hand tuned based on recommendations in [2])

$$\mathbf{PCEN}(t, f) = \left(\frac{\mathbf{E}(t, f)}{(\varepsilon + (\mathbf{E}^t * \phi_T)(t, f))^\alpha} + \delta \right)^r - \delta^r$$

Temporal integration

Automatic gain control

Dynamic range compression

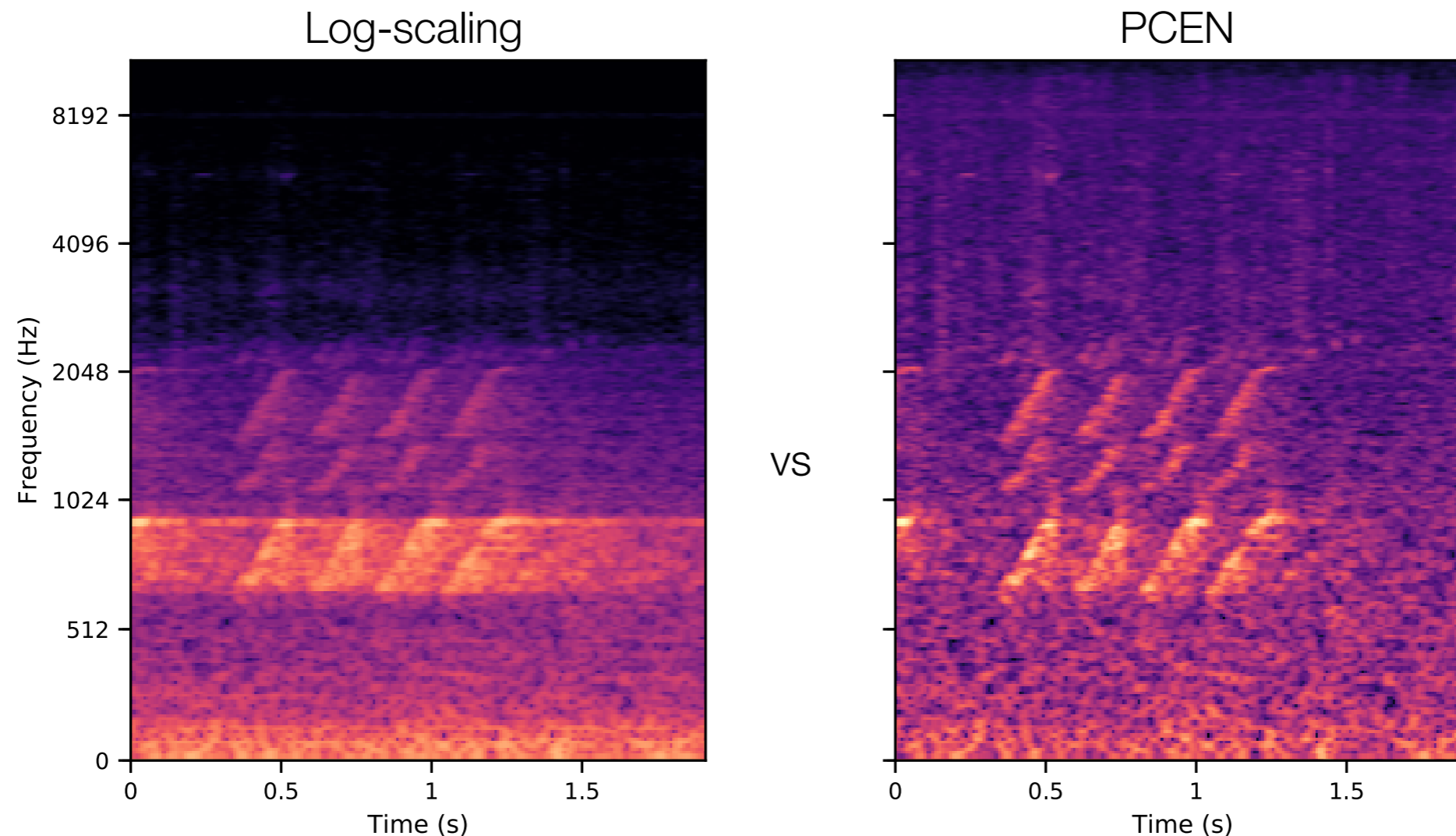
1. Wang, et al. "Trainable frontend for robust and far-field keyword spotting", ICASSP 2017

2. Lostanlen, Salamon, Cartwright, McFee, Farnsworth, Kelling, Bello, "Per-Channel Energy Normalization: Why and How", SPL 2019

Strategies to focus on foreground events: Per-Channel Energy Normalization (PCEN)

Pre-process with Per-Channel Energy Normalization (PCEN)¹

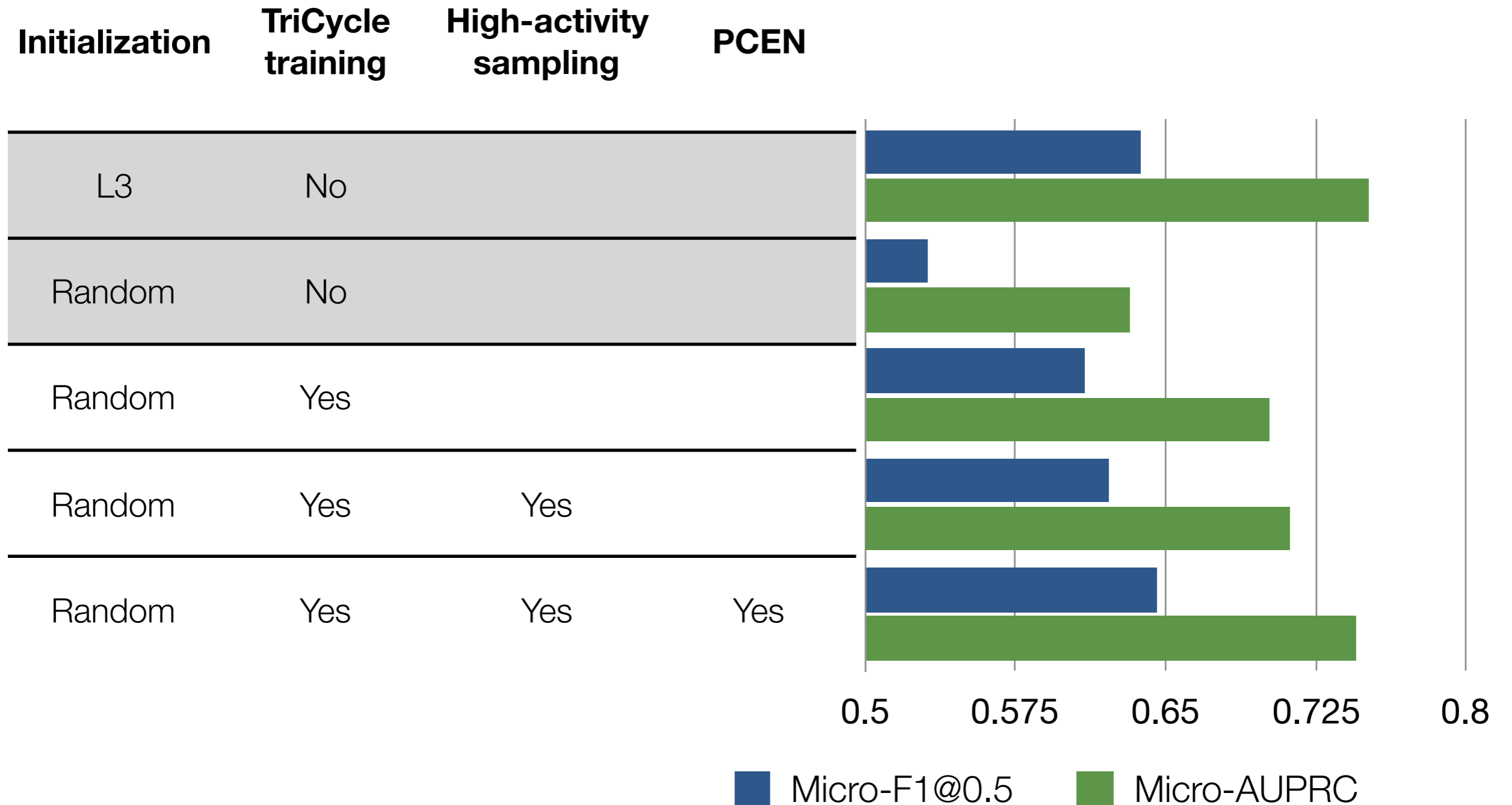
- Spectrogram processing that Gaussianizes and decorrelates frequency channels while retaining sound events of interest (parameter hand tuned based on recommendations in [2])



1. Wang, et al. "Trainable frontend for robust and far-field keyword spotting", ICASSP 2017

2. Lostanlen, Salamon, Cartwright, McFee, Farnsworth, Kelling, Bello, "Per-Channel Energy Normalization: Why and How", SPL 2019

Urban sound tagging results with TriCycle



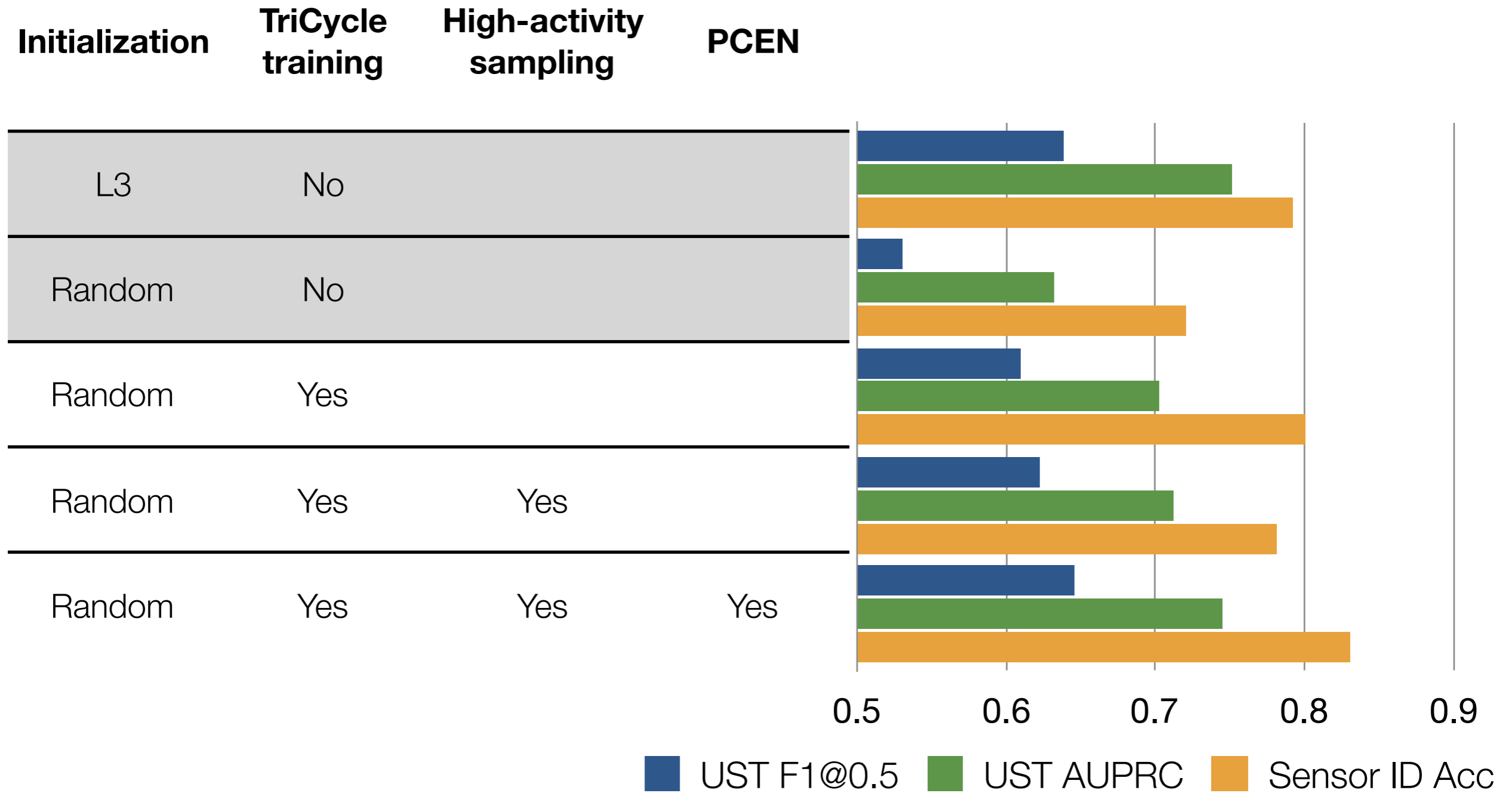
Future work

- Investigate circular regression loss formulations for von Mises distributed data
- Allow for groups of recordings with similar phase to be trained simultaneously and fused to increase the temporal signal and reduce impact of the background (hopefully reduce need for PCEN)
- Analyze the benefits of each temporal cycle and what information is encoded, and what is not
- Test TriCycle approach on other modalities

Summary

- Proposed an approach to self-supervised audio representation learning by predicting the time of recording
- First self-supervised embedding model trained on long-term temporal structure (regardless of modality)
- Able to train dataset-specific embeddings with single-modal data
- Validated approach on an urban sound tagging task, matching performance of a general state-of-the-art audio embedding
- Approach may be more general than audio, and well-suited for datasets from other sensor networks also having dense, longitudinal, timestamped data

Sensor prediction results with TriCycle



Results

Name	(a)			(b)			(c)				(d)
	Init.	TriCycle Train	Variation	MAD Day	MAD Week	MAD Year	UST F1@0.5	UST P@0.5	UST R@0.5	UST AUPRC	Sensor Acc.
<i>l3</i>	L ³ -Net	No	—	—	—	—	0.638	0.767	0.547	0.751	0.792
<i>rand</i>	Rand.	No	—	—	—	—	0.531	0.697	0.429	0.632	0.721
<i>rand-tc</i>	Rand.	Yes	—	0.480	0.508	0.562	0.622	0.734	0.540	0.712	0.781
<i>l3-tc-llr</i>	L ³ -Net	Yes	Low LR	0.370	0.531	0.540	0.638	0.764	0.548	0.739	0.824
<i>l3-tc-hlr</i>	L ³ -Net	Yes	High LR	0.338	0.443	0.545	0.638	0.749	0.556	0.737	0.851
<i>rand-tc-rs</i>	Rand.	Yes	Rand. Sampling	0.416	0.508	0.542	0.610	0.739	0.520	0.702	0.801
<i>rand-tc-pcen</i>	Rand.	Yes	PCEN	0.351	0.423	0.444	0.650	0.767	0.564	0.744	0.831