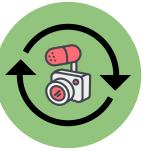
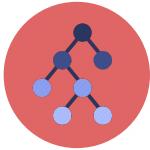


# Incorporating Natural Structure into Transfer Learning Methods for Machine Listening



## ECE PhD Dissertation Defense



Aurora Cramer  
(they/she)



Advised by  
Prof. Juan Pablo Bello, PhD



**NYU**

TANDON SCHOOL  
OF ENGINEERING

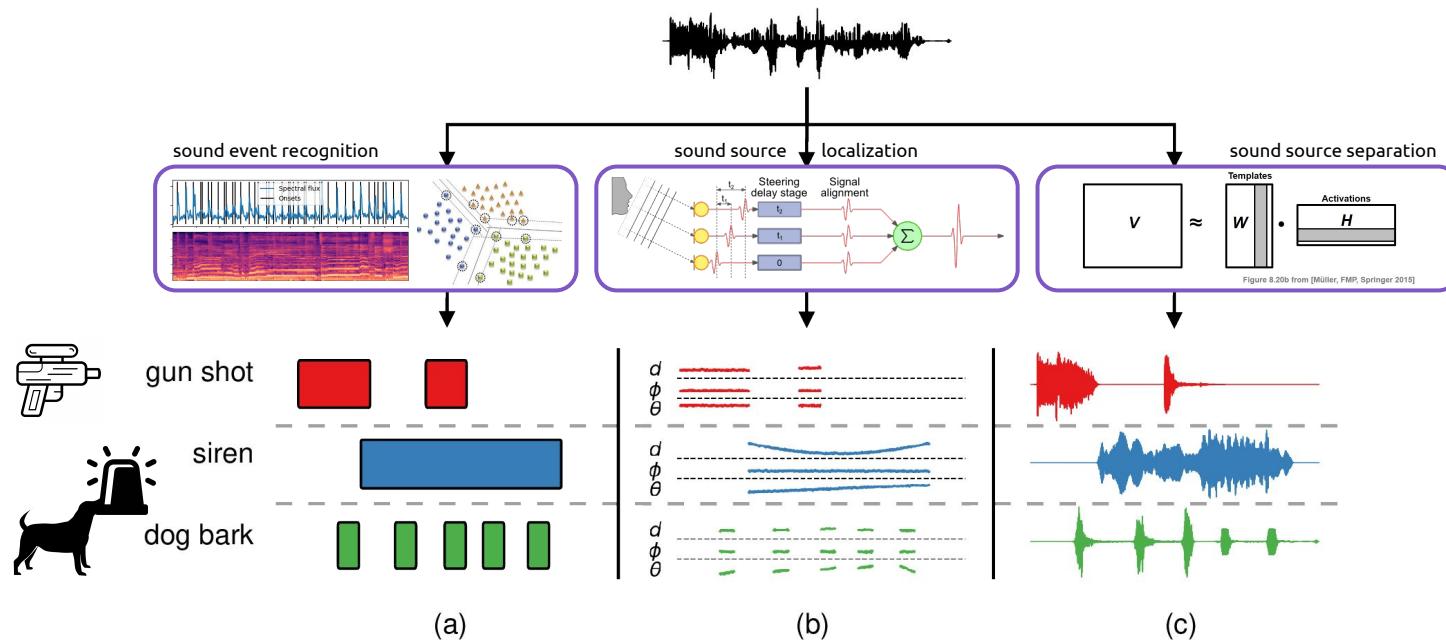


# Machine listening is everywhere!

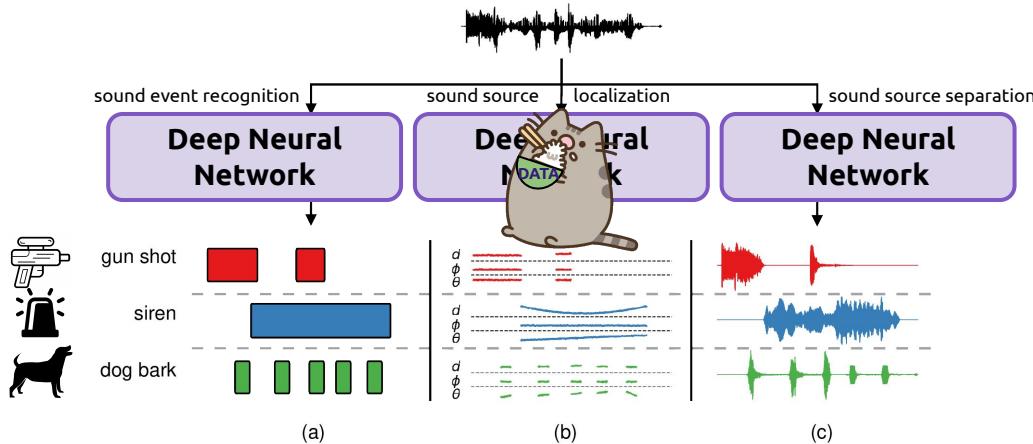


# Traditional machine listening methods

Traditional methods use handcrafted signal processing features and linear/Gaussian/Markov modeling methods



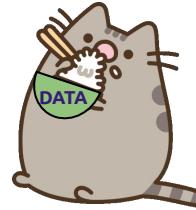
# Modern machine listening methods



Fully-supervised DNN models have generally pushed the SOTA for machine listening, but suffer from lack of abundant annotated audio data

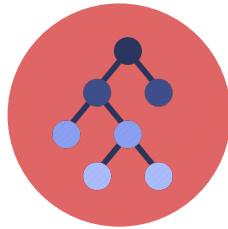


Audio annotation is time-consuming, expensive, and often difficult



How do we make more effective use of the data we have?

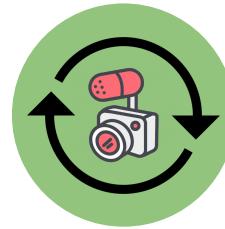
# What kinds of structure do people use?



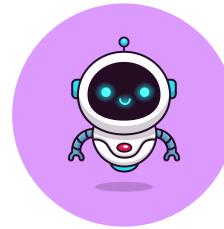
"It helps to know how concepts are related"



"similar tasks often require similar skills"



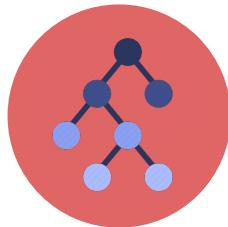
"many events can be perceived by multiple senses"



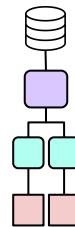
"learning can happen anywhere"

# What kinds of structure can machines use?

parallel transfer learning  
(a.k.a. multi-task learning)

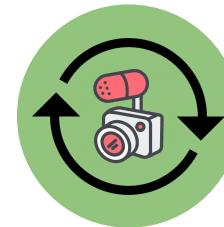


sound source  
hierarchies

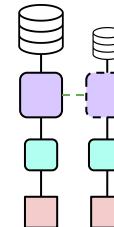


multi-task output  
structure

sequential transfer learning



multi-modal  
self-supervision

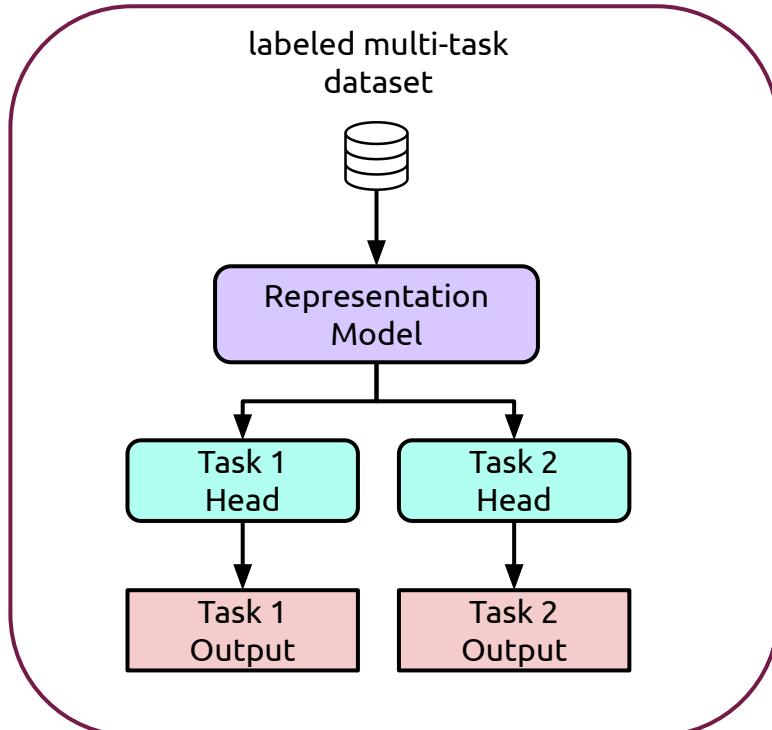


embodied agents

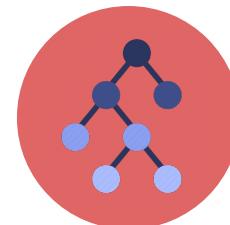
**Transfer Learning**

# Flavors of transfer learning

parallel transfer learning  
(a.k.a. multi-task learning)



- **Shared representation is regularized** by encoding relevant information for each task
- Model can learn **more efficiently** from the **same input data**



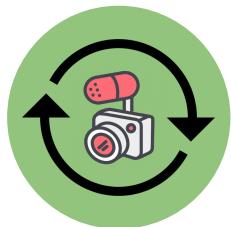
sound source  
hierarchies



multi-task output  
structure

# Flavors of transfer learning

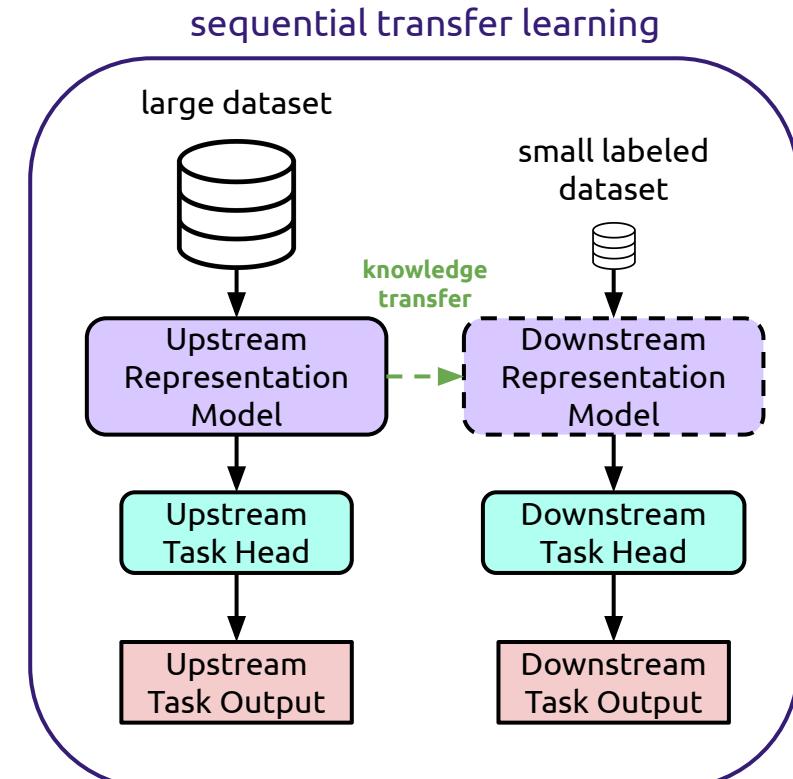
- Models for niche applications can **take advantage of existing large datasets** despite small target datasets
- Sufficiently **general representation** can be used to easily bootstrap downstream models for variety of domains



multi-modal  
self-supervision

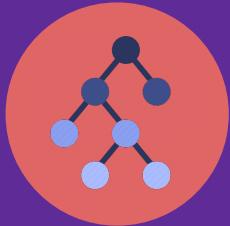


embodied agents



# Parallel transfer learning using sound source hierarchies

parallel transfer learning  
(a.k.a. multi-task learning)



sound source  
hierarchies



multi-task output  
structure

sequential transfer learning

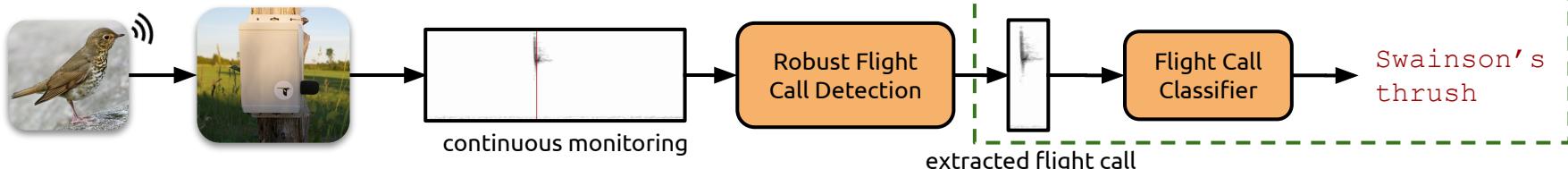


multi-modal  
self-supervision



embodied agents

# Automated flight call classification

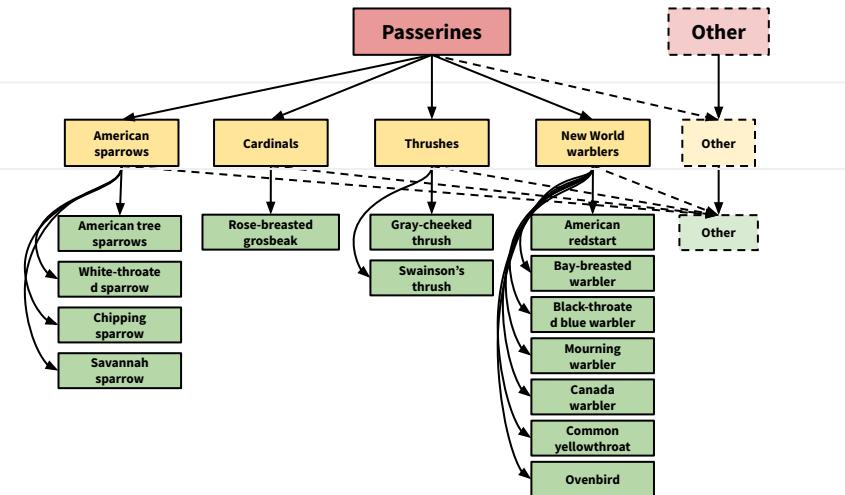
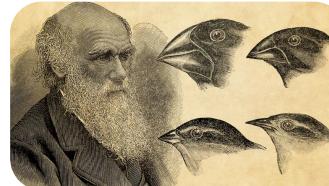


Incorporate hierarchical structure from **biological taxonomies** via **multi-task training** and **hierarchical model architectures**

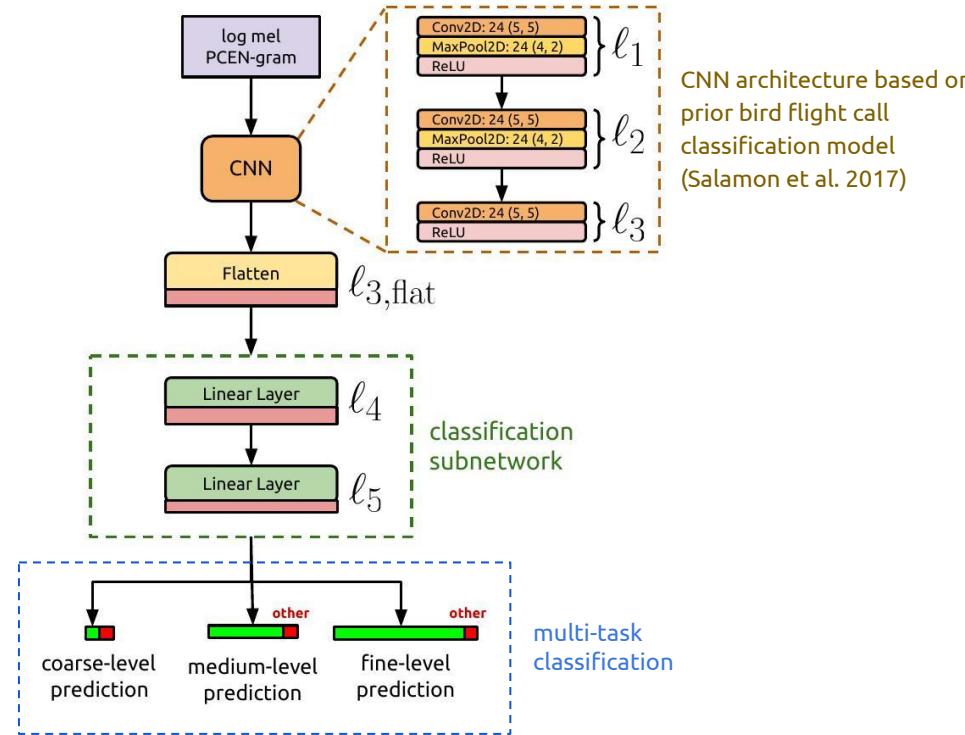
## Coarse Level: Order

## Medium Level: Family

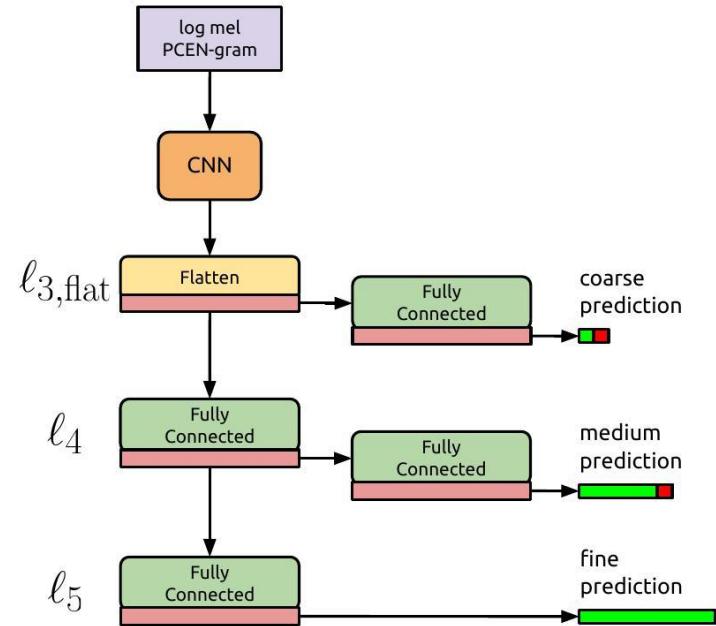
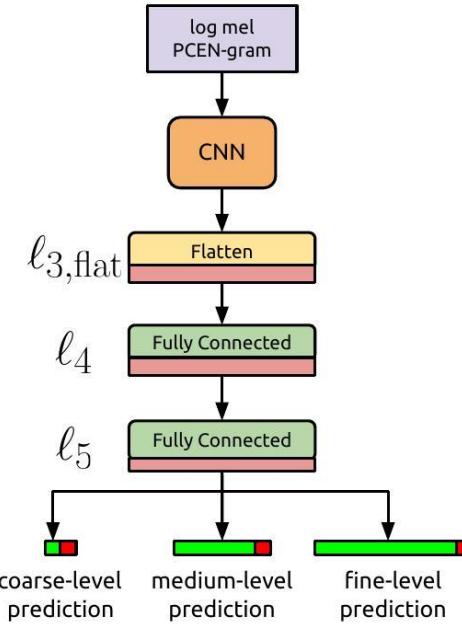
## Fine Level: Species



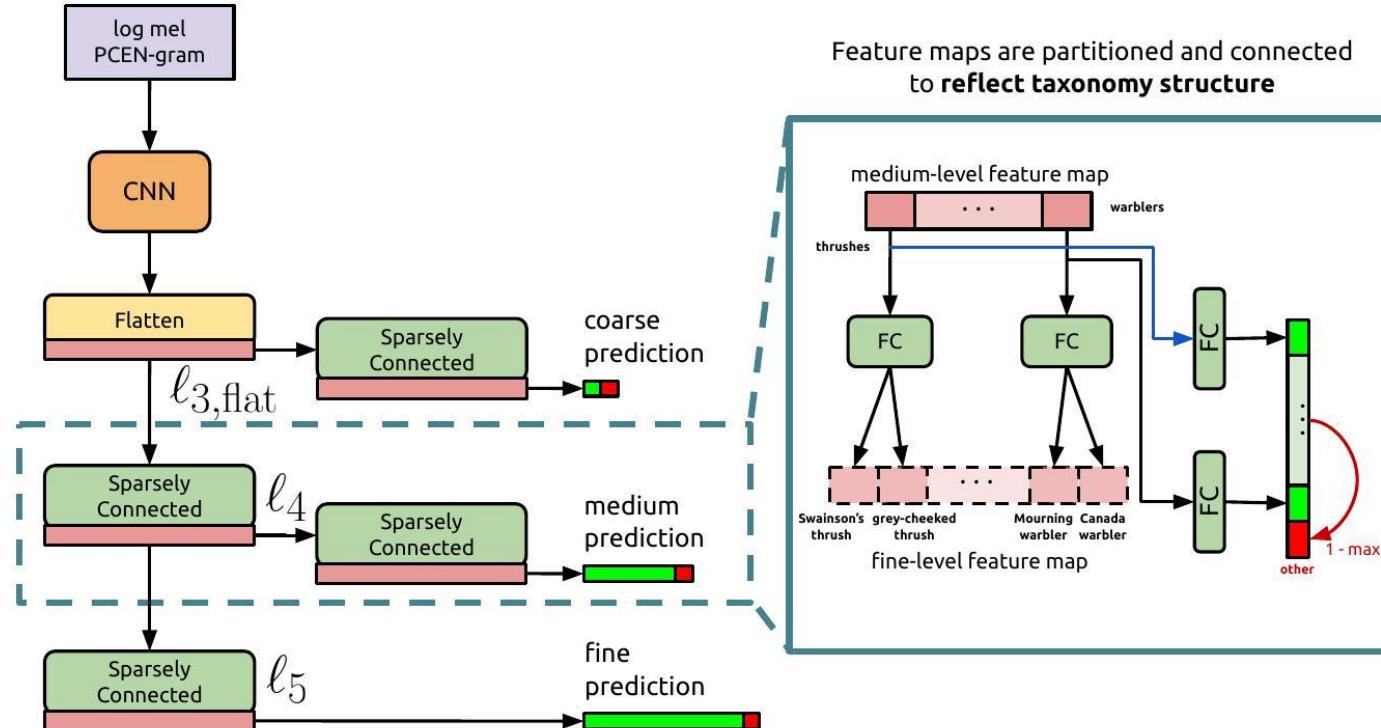
# Hierarchical flight call classification



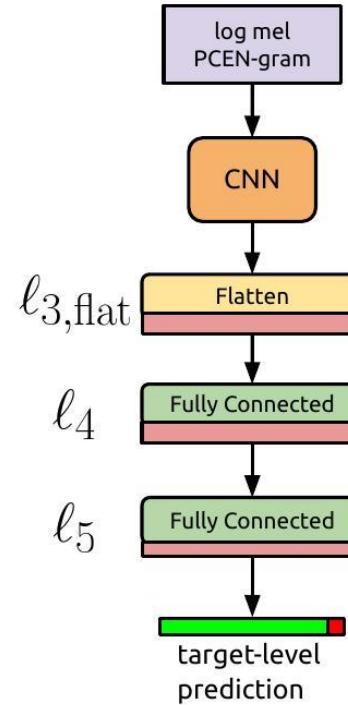
# Multi-task prediction (d) and hierarchical outputs (e)



# Hierarchical partitioning (f)

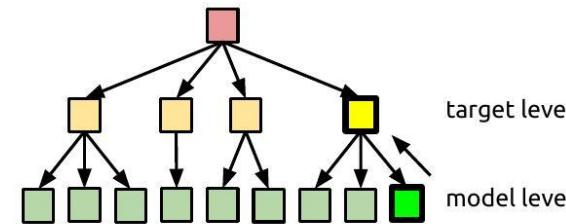


# Single-task baseline (a) - (c)



**Specialist strategy:** use model trained at target level

**Coarsening strategy:** use target-level ancestor of predicted taxa



# Training and evaluation

- Train on **heterogeneous dataset** of flight calls from various sources (ANAFCC<sup>★</sup>)
- Input to model: 150 ms **log-scale mel-frequency spectrogram** with **per-channel energy normalization (PCEN)** applied
- **Data augmentation:** pitch shifting, time-stretching, additive background noise
- Evaluate at each taxonomic level using annotated clips from **full season of sensor network recordings** (BirdVox-14SD<sup>‡</sup>)

★ ANAFCC: <https://doi.org/10.5281/zenodo.3666782>

‡ BirdVox-14SD: <https://doi.org/10.5281/zenodo.3667094>

# Experimental Results

Multi-task models  
match specialized  
and outperform  
coarsened!

safe, simple choice

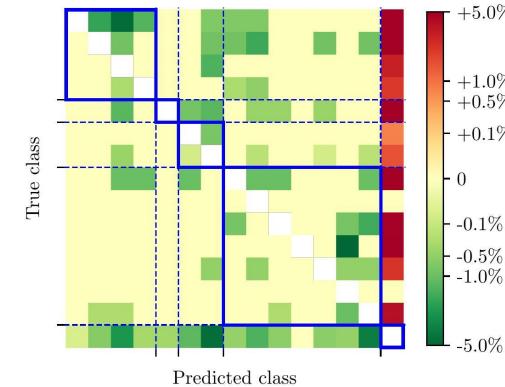
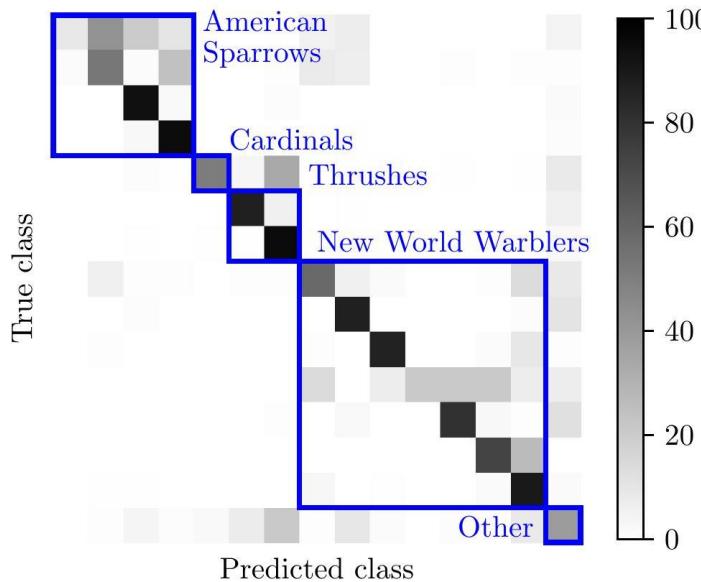
better for variety of species

better for common species

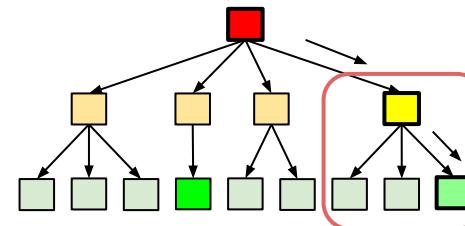
better for coarser taxa

Model					# Trained Params	Fine Micro Acc.	Fine Macro Acc.	Medium Micro Acc.	Medium Macro Acc.	Coarse Acc.
<u>Single-Task Model</u>										
(a)	Fine Level				641K	61.13	54.80	64.61	50.40	77.72
(b)	Medium Level				640K	-	-	73.80	56.04	<b>94.75</b>
(c)	Coarse Level				640K	-	-	-	-	93.85
<u>TaxoNet Model</u>										
	Layer Partitioning	Hier. Outputs	Classifier Activation	Classifier Projection						
(d)	No	No	sigmoid	Trainable	641K	61.82	55.83	75.10	55.87	94.39
(e)	No	Yes	sigmoid	Trainable	650K	58.74	<b>58.06</b>	75.83	60.04	94.54
(f)	Yes	Yes	sigmoid	Trainable	649K	<b>66.33</b>	55.69	76.50	61.60	94.69
(g)	Yes	Yes	softmax	Trainable	649K	60.39	52.30	75.94	56.96	94.67
(h)	Yes	Yes	tanh	Mean	640K	63.47	41.46	<b>79.36</b>	<b>65.08</b>	<b>94.75</b>

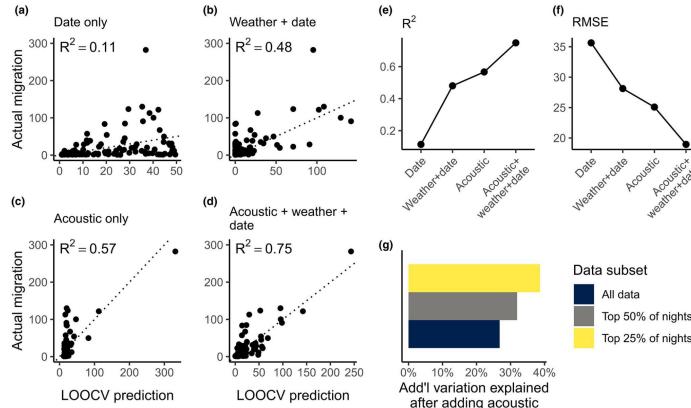
# Model makes reasonable errors



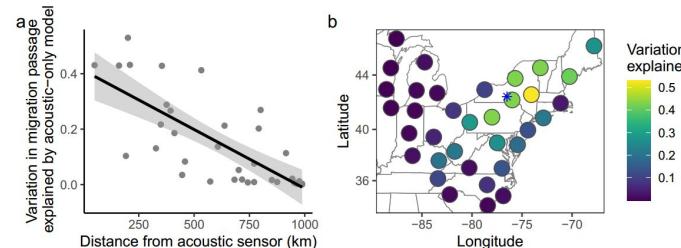
Enforcing hierarchical consistency procedure reduces confusion further



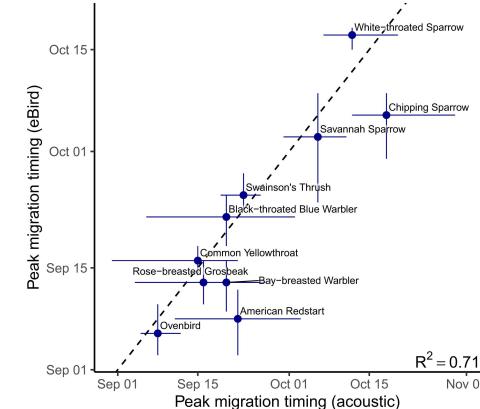
# Automated acoustic monitoring captures intensity and timing of bird migration!



generalized linear model predicts migration intensity as measured by weather radar



model also captures regional intensity at nearby radar locations



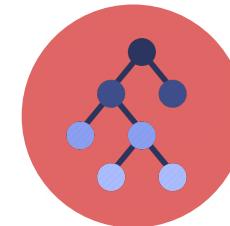
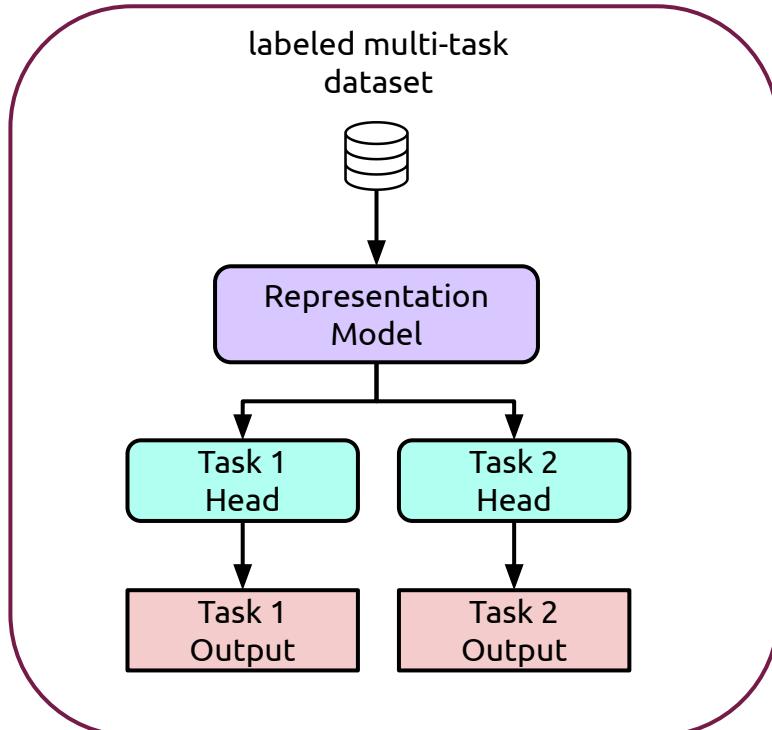
generalized additive model predictions of peak migration time correlate with birdwatcher estimates



```
$ pip install birdvoxclassify
$ pip install birdvoxdetect
```

# Flavors of transfer learning

parallel transfer learning  
(a.k.a. multi-task learning)

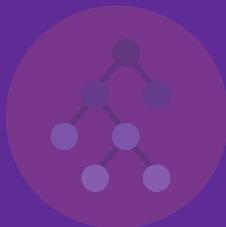


sound source  
hierarchies

**hierarchical structure**  
improves **model robustness!**

# Parallel transfer learning using multi-task output structure

parallel transfer learning  
(a.k.a. multi-task learning)



sound source  
hierarchies



multi-task output  
structure

sequential transfer learning

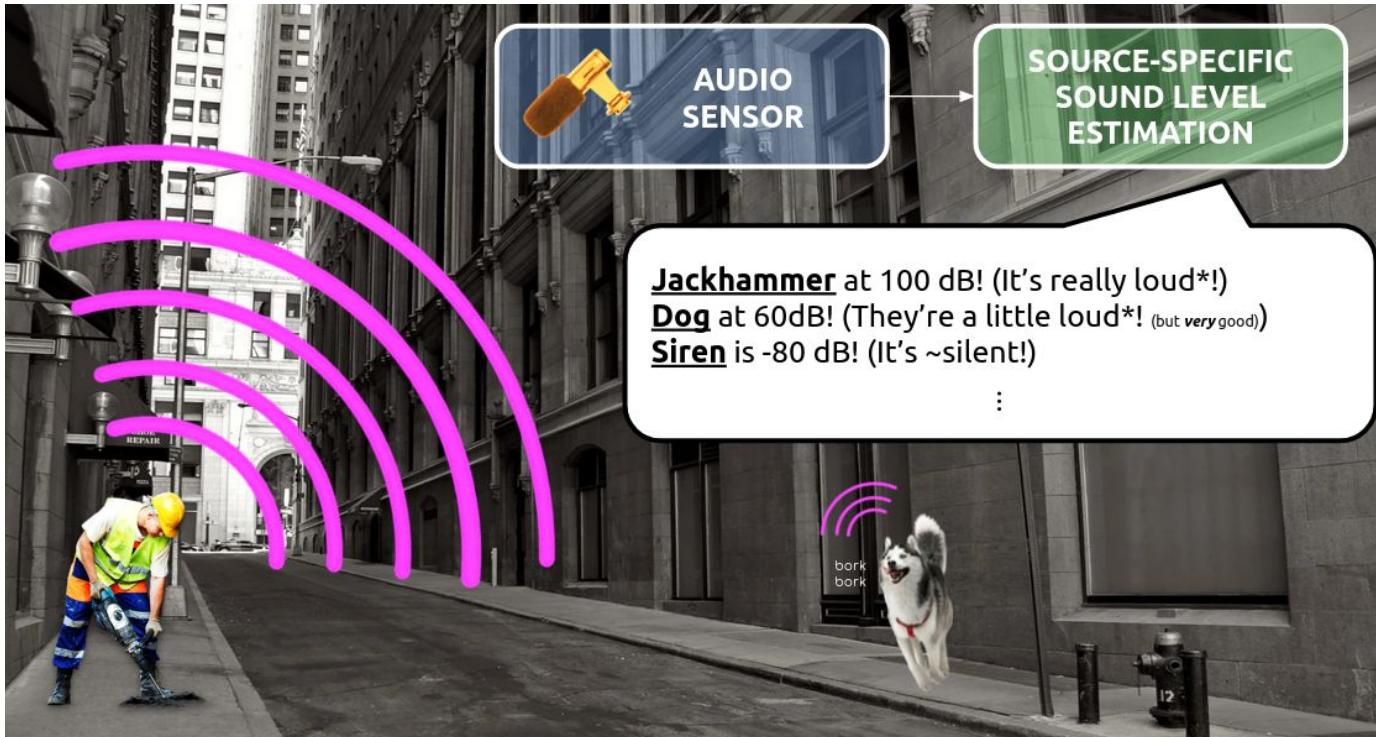


multi-modal  
self-supervision

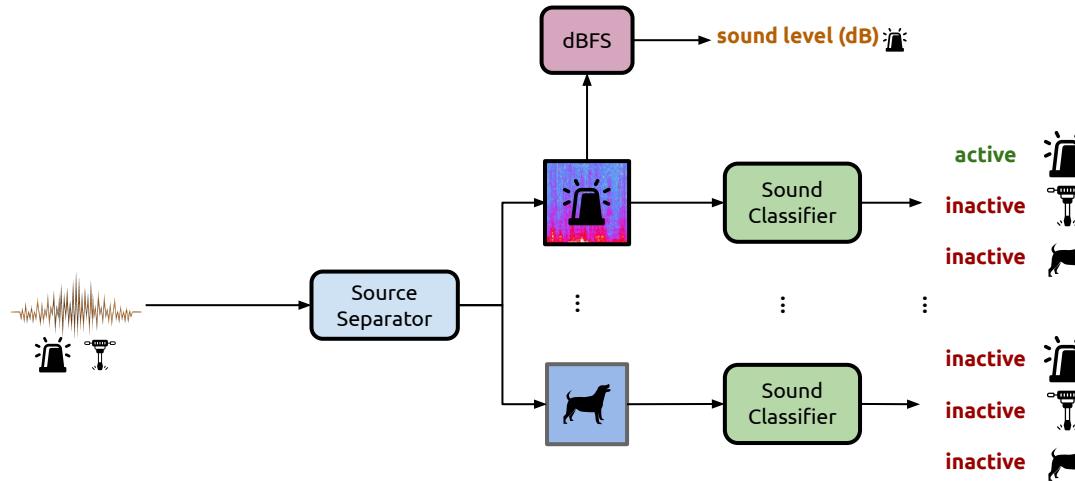


embodied agents

# Source specific sound level estimation is difficult in practice

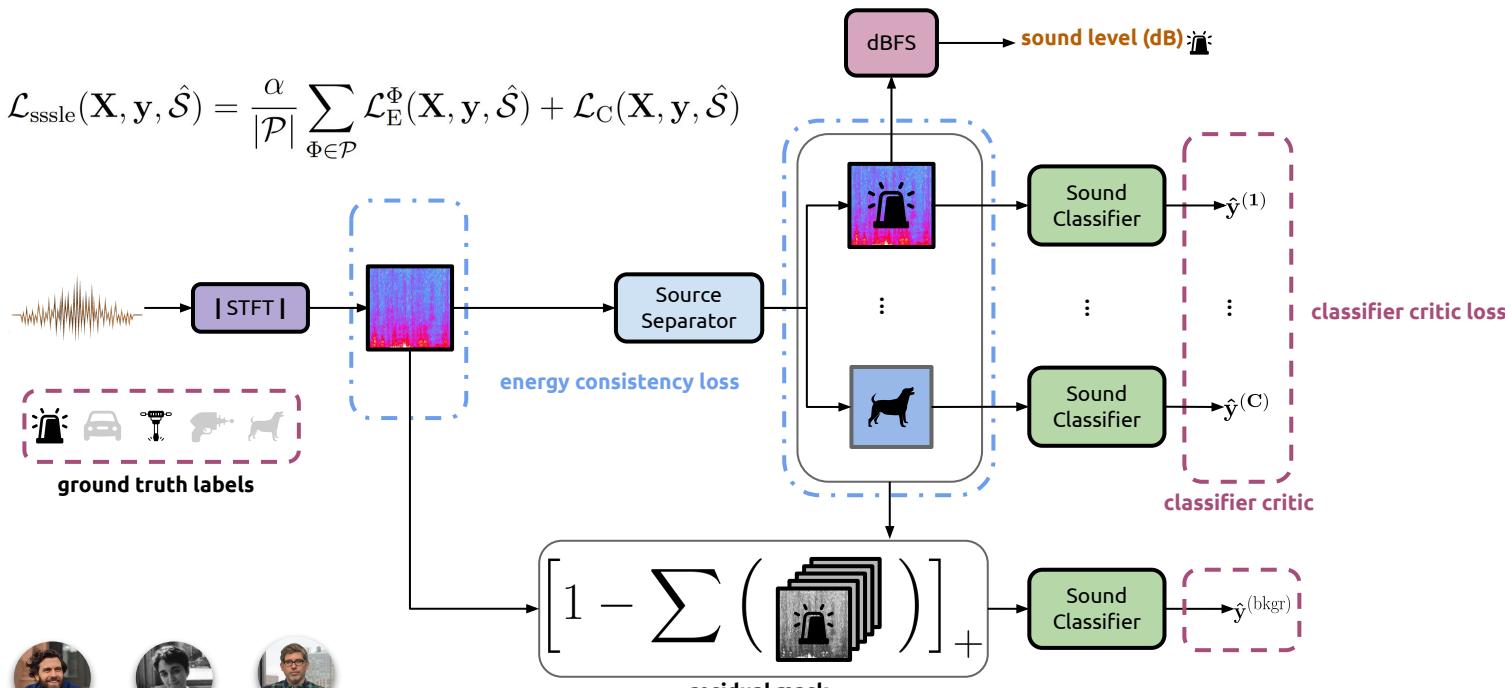


# Source specific sound level estimation with multi-task learning

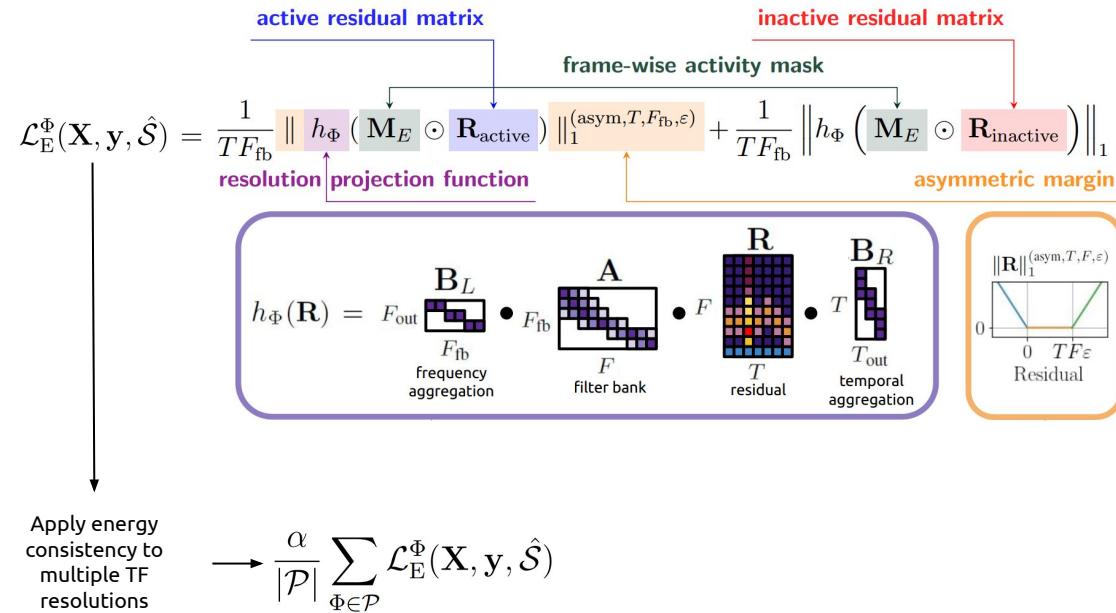
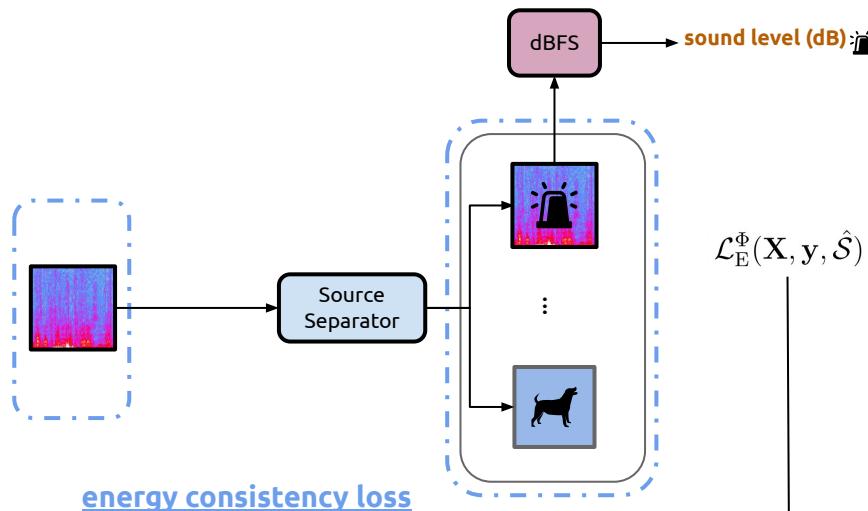


# Source specific sound level estimation with multi-task learning

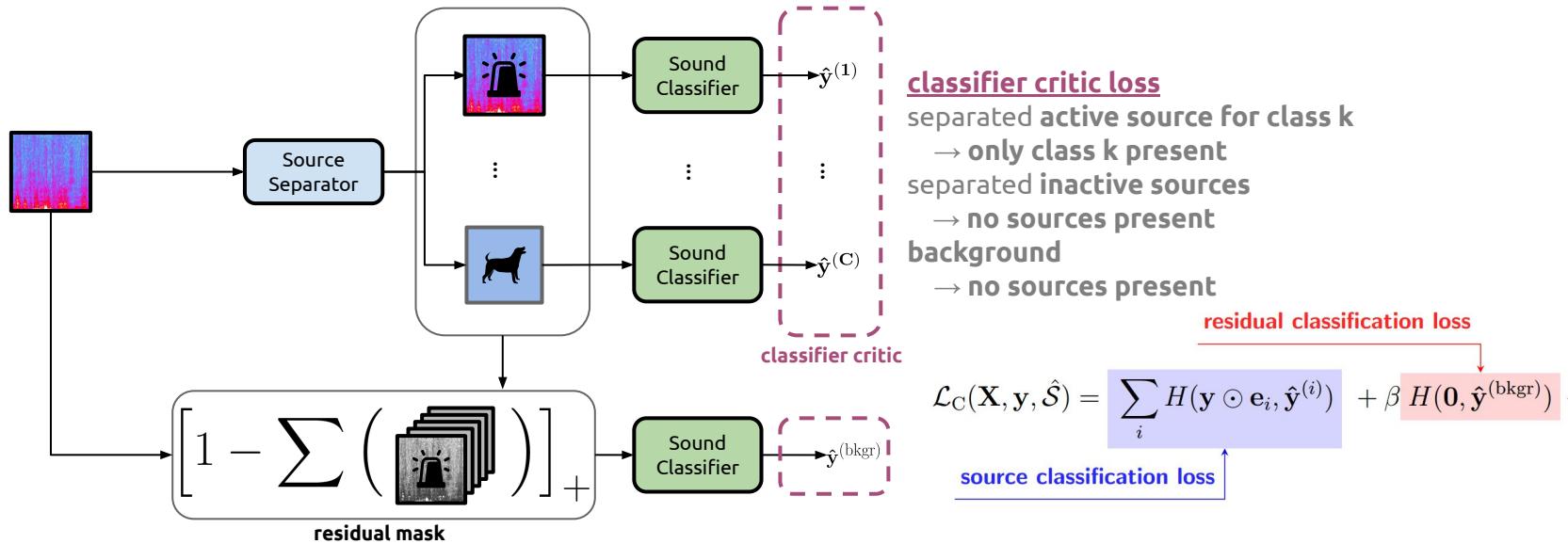
$$\mathcal{L}_{\text{ssse}}(\mathbf{X}, \mathbf{y}, \hat{\mathcal{S}}) = \frac{\alpha}{|\mathcal{P}|} \sum_{\Phi \in \mathcal{P}} \mathcal{L}_E^\Phi(\mathbf{X}, \mathbf{y}, \hat{\mathcal{S}}) + \mathcal{L}_C(\mathbf{X}, \mathbf{y}, \hat{\mathcal{S}})$$

NYU  
TANDON SCHOOL  
OF ENGINEERINGMARL  
New Jersey Institute  
of TechnologyNJIT  
Northwestern  
University

# Source specific sound level estimation with multi-task learning



# Source specific sound level estimation with multi-task learning



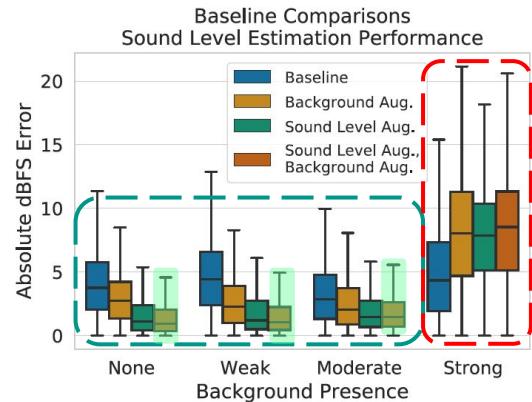
## Training and evaluation

- Train and evaluate with 4-second synthetic soundscapes (16kHz) containing **urban sound events** (from UrbanSound8K) and different levels of **urban background noise** (SONYC-Backgrounds \*)
  - Different background levels: -50/-20/0 dB LUFS (weak/moderate/strong), no background
- Evaluate with respect to both **source separation performance** (SI-SDR improvement) and **sound level estimation** performance (absolute dBFS error), comparing with weakly supervised source separation baseline

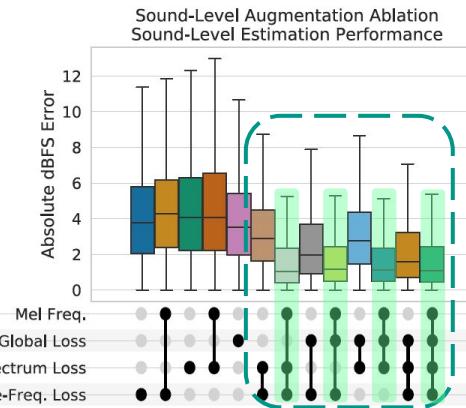
\* SONYC-Backgrounds: <https://doi.org/10.5281/zenodo.5129078>

Salamon et al., "A dataset and taxonomy for urban sound research," ACM Multimedia, 2014.

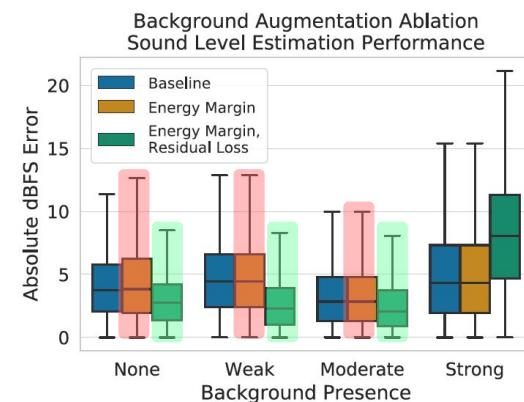
# Experimental Results



**SSSLE performance (and source separation performance) are improved in up to moderate background conditions!**



**Using multiple time-frequency resolutions is beneficial**

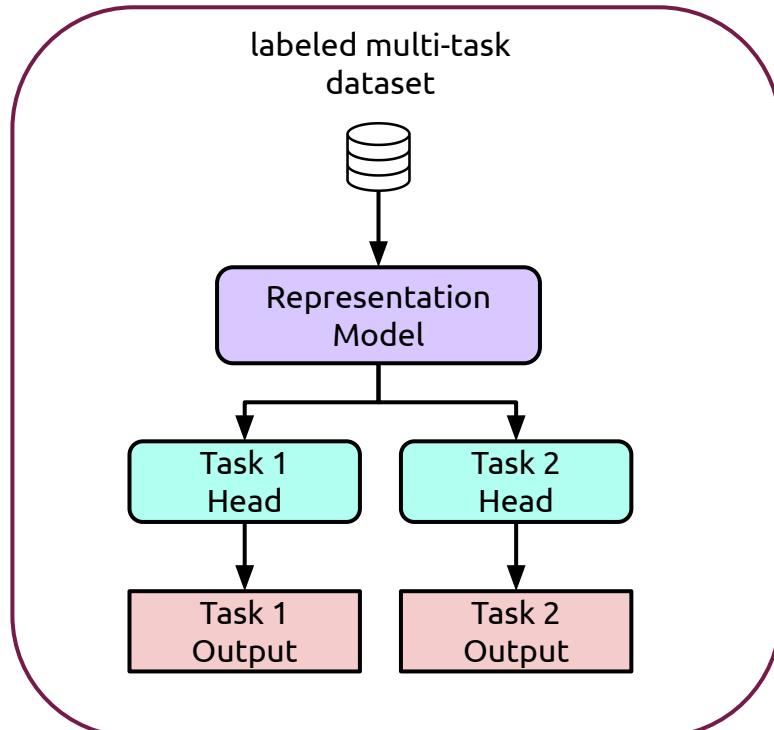


**The background classification loss is crucial when using an asymmetric margin**

**SSSLE models can be trained with only clip-level source labels using multi-task learning!**

# Flavors of transfer learning

parallel transfer learning  
(a.k.a. multi-task learning)



multi-task output  
structure

identifying and capitalizing on  
**multi-task output structure**  
improves **cross-task transfer!**

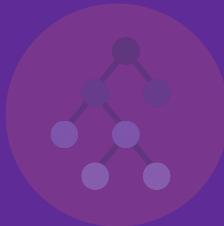
# What if there still isn't enough data?

We can transfer knowledge from **other models** that have already learned useful representations!



# Sequential transfer using multi-modal self-supervision

parallel transfer learning  
(a.k.a. multi-task learning)

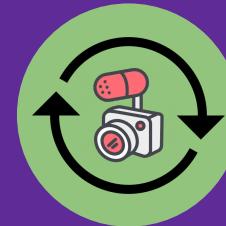


sound source  
hierarchies



multi-task output  
structure

sequential transfer learning

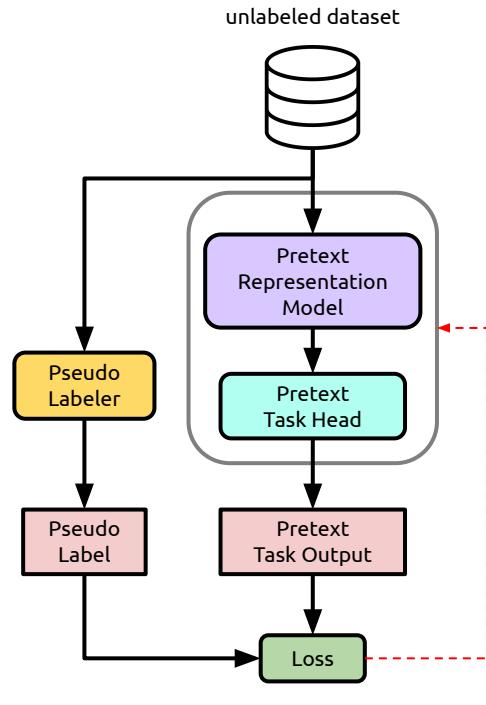


multi-modal  
self-supervision

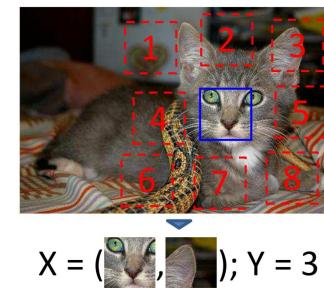


embodied agents

# Self-supervised transfer learning

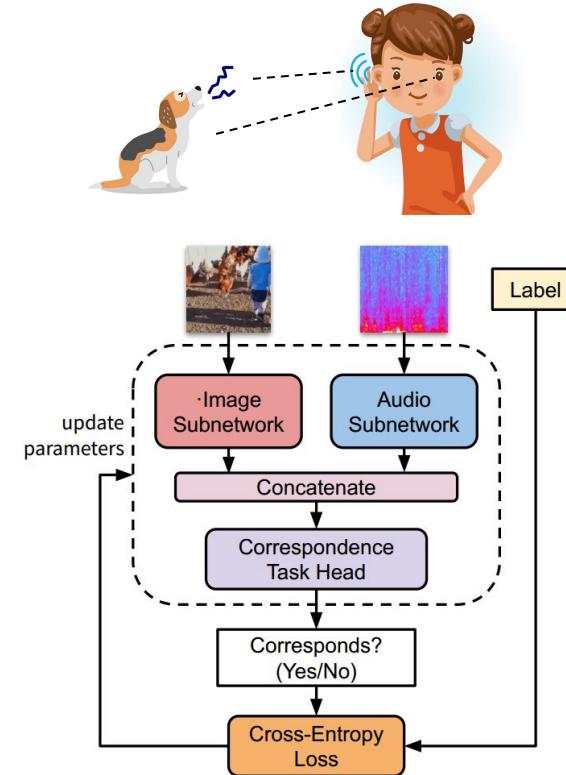


- Train models with **unlabeled datasets** using “**pretext**” tasks with corresponding **pseudo-labels**
- Model **implicitly learns representation** encoding useful perceptual and/or semantic information

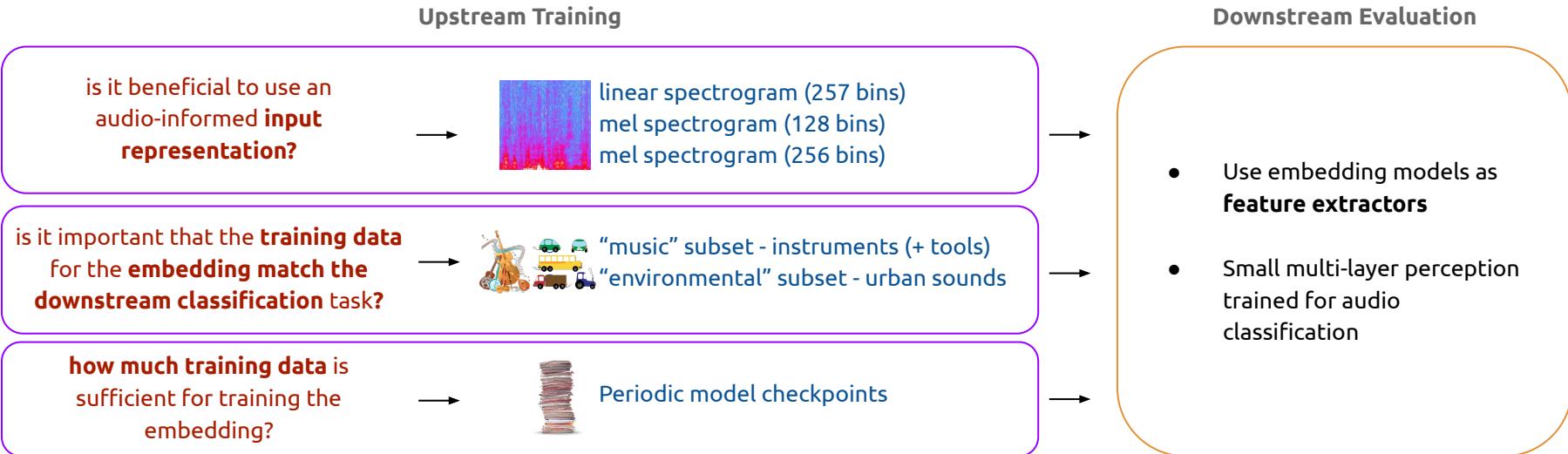


# Audio-visual correspondence (AVC)

- Auditory and visual stimuli often correspond to a common underlying source → *audio-visual correspondence*
- AVC defines a **simple self-supervision task** that learns embeddings from **unlabeled videos** that are useful for both **downstream audio and visual tasks!**
- Can we better understand how to train effective audio embeddings using AVC?



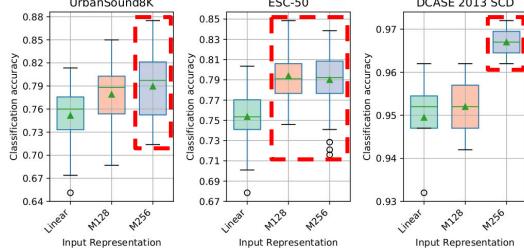
# How do design choices affect downstream audio classification performance?



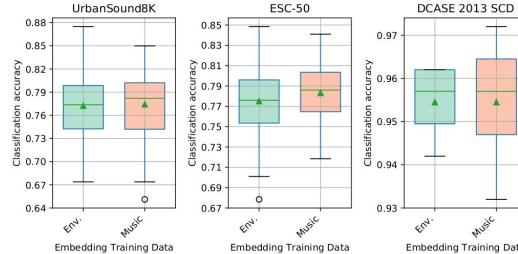
# Training and evaluation

- Train and validate upstream AVC model using subsets of AudioSet
  - Generate pairs of 1 second audio clips and random overlapping video frame for positive, shuffle pairs for negatives
  - Data augmentation - video: crop, brightness, contrast, saturation; audio: loudness)
- Train and evaluate downstream models using multi-class audio classification datasets (UrbanSound8K, ESC-50, DCASE 2013 Scene Classification)
  - Obtain clip-level predictions by averaging framewise predictions on 1 second overlapping frames

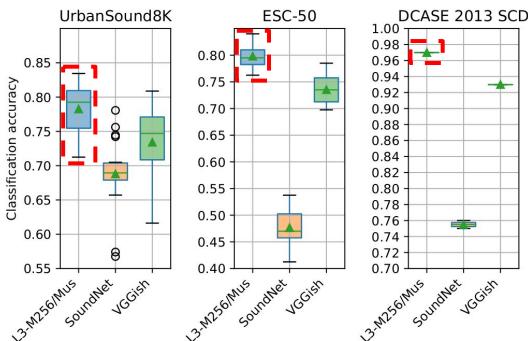
# Experimental Results



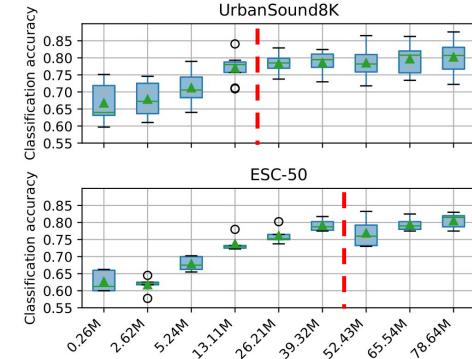
Mel spectrograms with half the number of bins outperform linear spectrograms!



Training with either subsets produces similar results — semantic content is less relevant than a strong AVC signal



L3-Net outperforms other competitive embeddings!

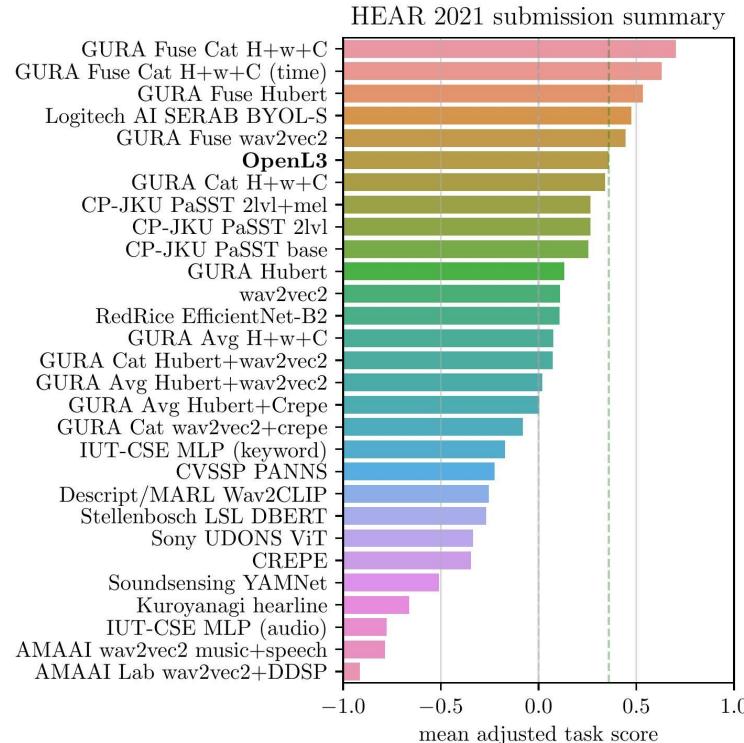


25-50% of AVC examples needed to match best performance



\$ pip install openl3

# OpenL3 remains competitive!



# Urban sound exhibits natural urban rhythm



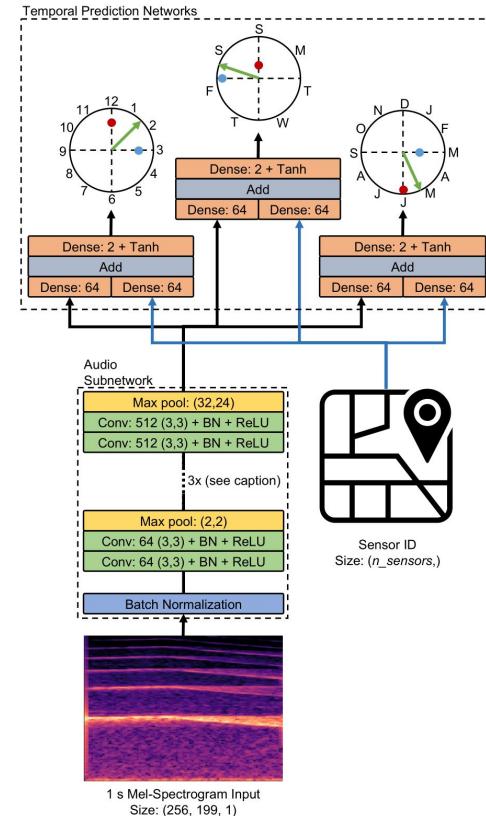
Soundscapes often sound different depending on the **time of day**,  
**day of the week**, or **month of the year**

# Temporal cycle prediction

- Leverage **self-supervised learning** with an L3-like architecture using ***temporal cycle prediction*** as the **pretext task**
- Model predicts the phase in **daily**, **weekly**, and **yearly cycles** using only sensor network audio
- Condition on sensor ID to help prevent model from overfitting to sensor characteristics



M. Cartwright, A. Cramer, J. Salamon, and J.P. Bello. "TriCycle: Audio representation learning from sensor network data using self-supervision," WASPAA, 2019



# Training and evaluation

- Train with **unlabeled clips** from a year of **urban soundscape recordings** obtained from the SONYC acoustic sensor network
  - Sample evenly in time, focusing on potentially meaningful events by randomly selecting recordings for each hour in the top 15th percentile of SPL difference:
$$\sqrt{\sum_{n=0}^{79} (d_{m,n} - d_{m,n-1})^2}$$
- Evaluate with a labeled **urban sound tagging dataset** from temporally-disjoint SONYC recordings (SONYC-UST v1)

M. Cartwright, A. Cramer, J. Salamon, and J.P. Bello. "TriCycle: Audio representation learning from sensor network data using self-supervision," WASPAA, 2019

M. Cartwright, A. E. M. Mendez, G. Dove, A. Cramer, V. Lostanlen, H. Wu, J. Salamon, O. Nov, and J. P. Bello, "SONYC urban sound tagging (SONYC-UST): a multilabel dataset from an urban acoustic sensor network," Zenodo, (2019). <https://zenodo.org/record/2590742>

# Experimental Results

Able to produce embeddings **comparable with L3-Net**  
on downstream urban sound tagging

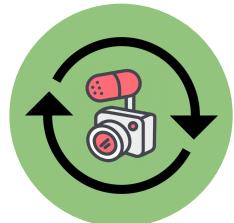
Name	Init.	TriCycle Train	Variation	MAD Day	MAD Week	MAD Year	UST F1@0.5	UST P@0.5	UST R@0.5	UST AUPRC	Sensor Acc.
<i>l3</i>	L <sup>3</sup> -Net	No	—	—	—	—	0.638	<b>0.767</b>	0.547	<b>0.751</b>	0.792
<i>rand</i>	Rand.	No	—	—	—	—	0.531	0.697	0.429	0.632	0.721
<i>rand-tc</i>	Rand.	Yes	—	0.480	0.508	0.562	0.622	0.734	0.540	0.712	0.781
<i>l3-tc-llr</i>	L <sup>3</sup> -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 <sup>3</sup> -Net	Yes	High LR	<b>0.338</b>	0.443	0.545	0.638	0.749	0.556	0.737	<b>0.851</b>
<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	<b>0.423</b>	<b>0.444</b>	<b>0.650</b>	<b>0.767</b>	<b>0.564</b>	0.744	0.831



M. Cartwright, A. Cramer, J. Salamon, and J.P. Bello. "TriCycle: Audio representation learning from sensor network data using self-supervision," WASPAA, 2019

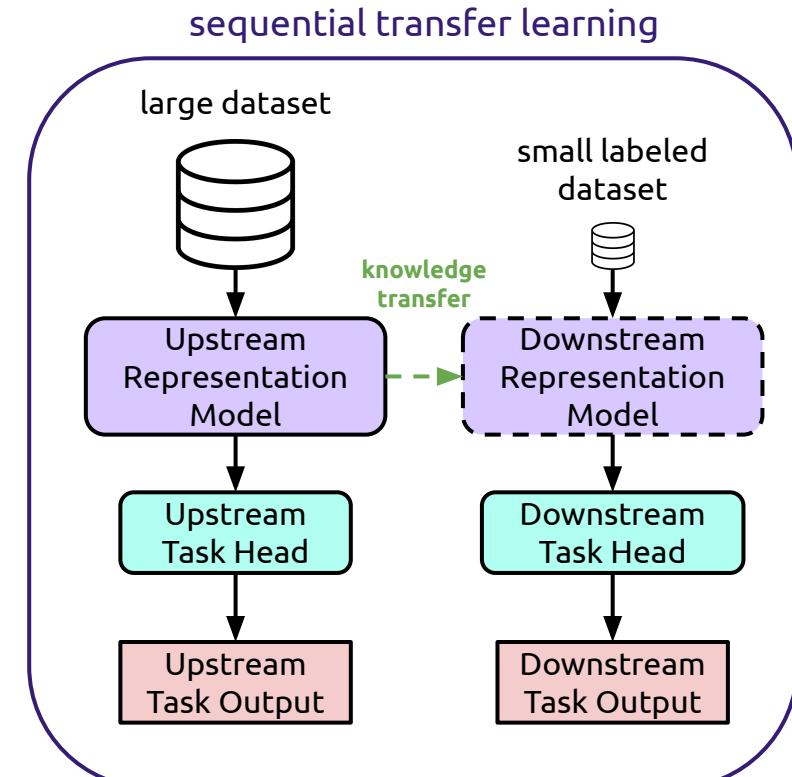
M. Cartwright, A. E. M. Mendez, G. Dove, A. Cramer, V. Lostanlen, H. Wu, J. Salamon, O. Nov, and J. P. Bello, "SONYC urban sound tagging (SONYC-UST): a multilabel dataset from an urban acoustic sensor network," Zenodo, (2019). <https://zenodo.org/record/2590742>

# Flavors of transfer learning

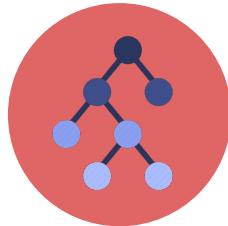


multi-modal  
self-supervision

natural **multi-modal correspondence** provide structure **without human annotations!**



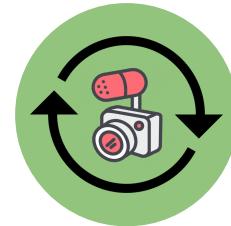
# What's next?



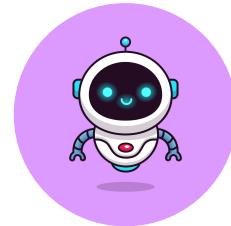
sound source  
hierarchies



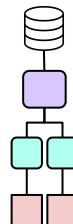
multi-task output  
structure



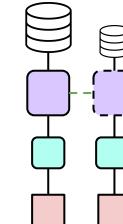
multi-modal  
self-supervision



embodied agents

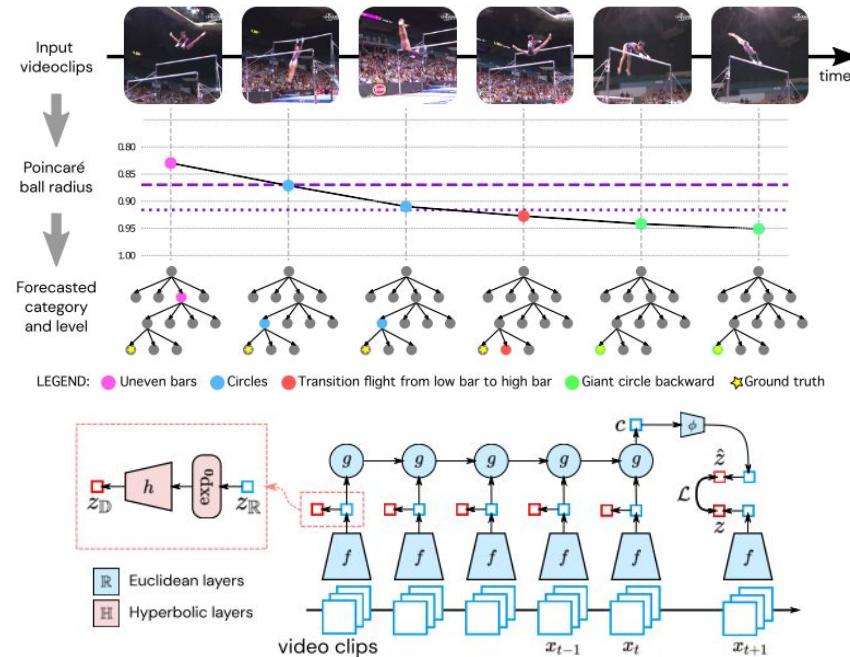


# Transfer Learning

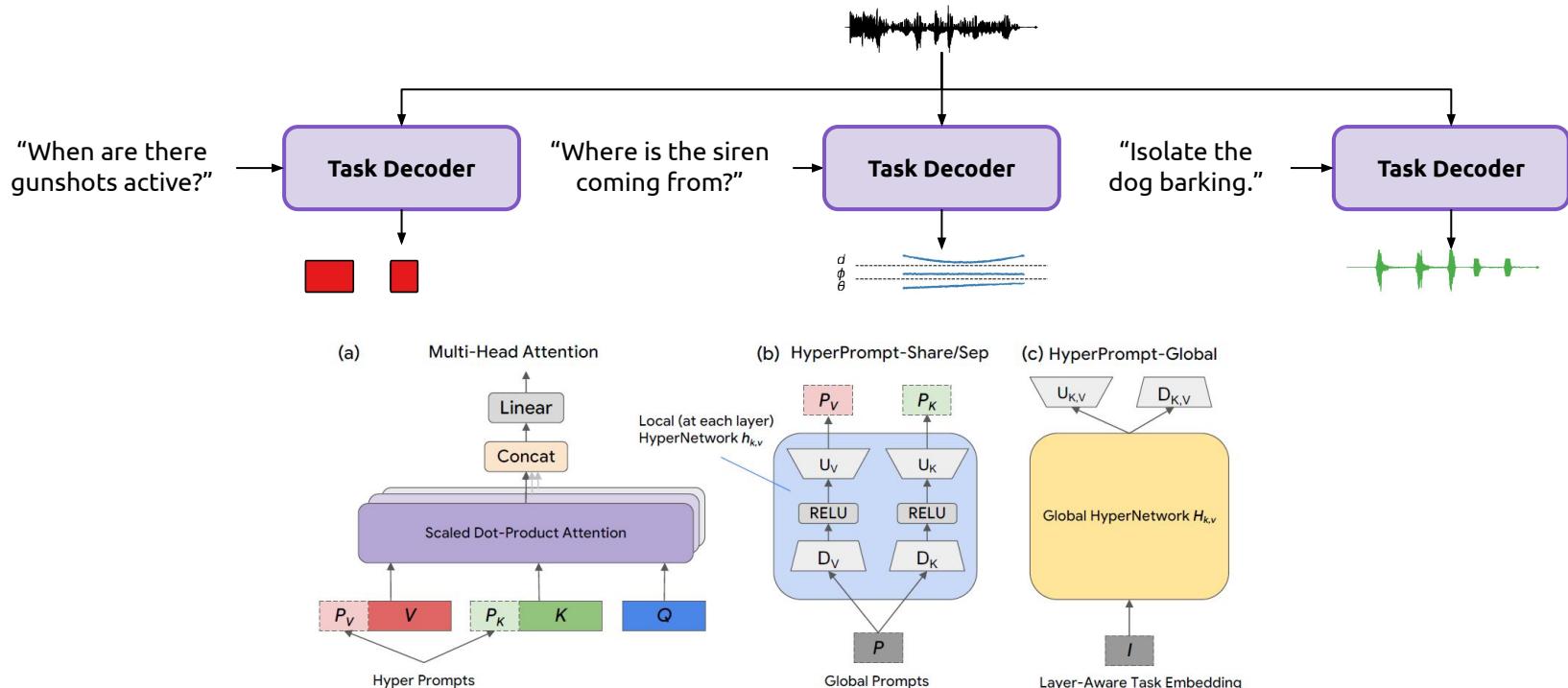


# Beyond predefined hierarchies

- Unlike Euclidean embeddings, **Hyperbolic embeddings** can directly encode hierarchical structure without distortion
- **Self-supervised contrastive predictive coding** with hyperbolic embeddings can implicitly learn hierarchies and uncertainty



# Beyond predefined tasks



# Beyond audio-visual correspondence

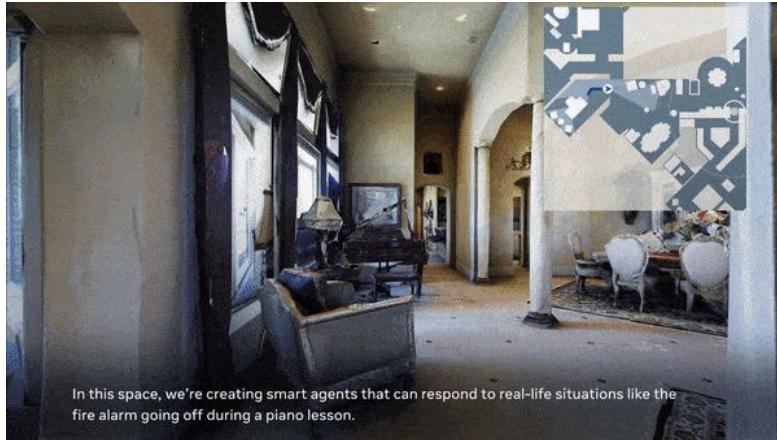
- Egocentric videos may better align to human perception than videos filmed with handheld devices
- Telemetry data like accelerometry, recording timestamps, and location can be used for further self-supervision
- Embodied agents that interact with the environment may even more effectively align to human perception



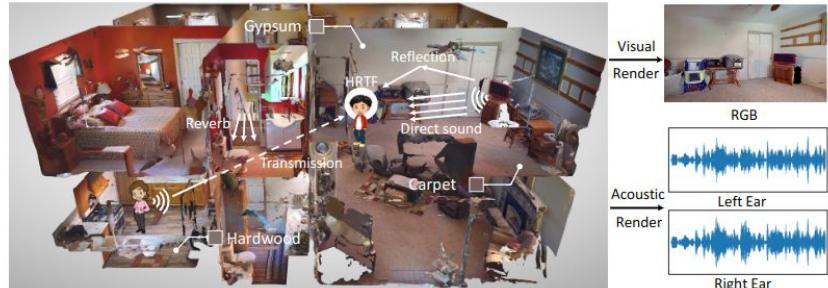
K. Grauman et al., "Ego4D: Around the World in 3,000 Hours of Egocentric Video," 2021.

S. Zhang et al., "EgoBody: Human Body Shape and Motion of Interacting People from Head-Mounted Devices," ECCV, 2022.

# Embodied navigation



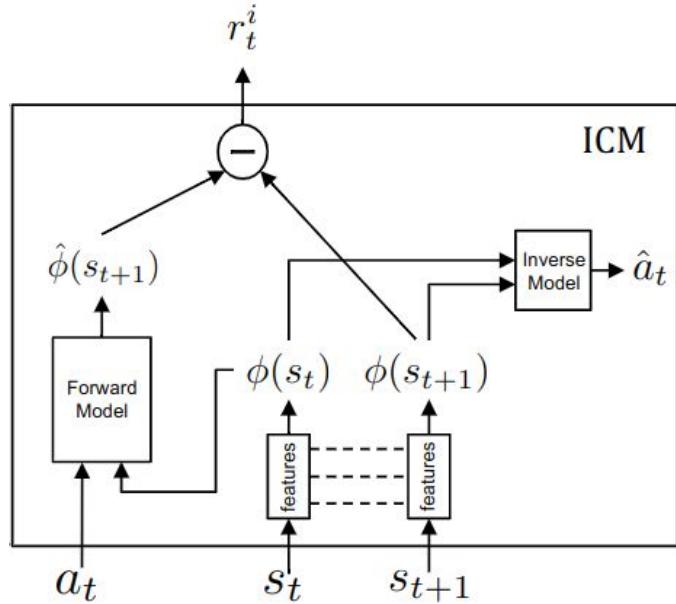
Photorealistic renders from point clouds  
+ spatialized audio from realistic RIRs



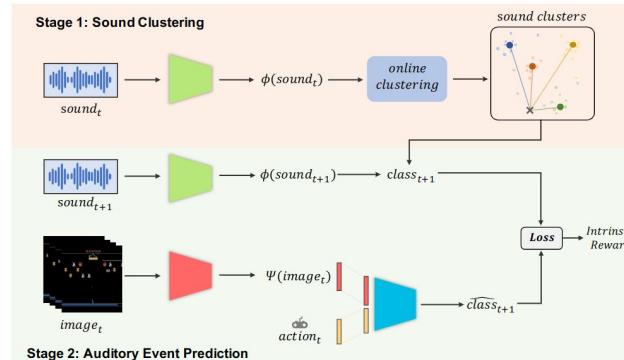
Agent must navigate to audible goal using audio and visual sensory inputs

- C. Chen *et al.*, “SoundSpaces 2.0: A Simulation Platform for Visual-Acoustic Learning,” *ArXiV*, 2022  
C. Chen *et al.*, “SoundSpaces: Audio-Visual Navigation in 3D Environments,” 2020.  
C. Gan *et al.*, “ThreeDWorld: A Platform for Interactive Multi-Modal Physical Simulation,” *NeurIPS*, 2021

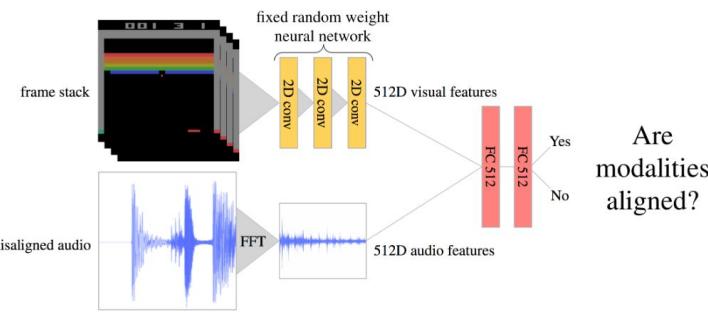
# Self-supervision through intrinsic valuation



Intrinsic Curiosity Module: predict action from state embeddings before and after action



Predict auditory events from visual stream



Are modalities aligned?

D. Pathak *et al.*, "Curiosity-driven Exploration by Self-supervised Prediction," *ICML*, 2017.

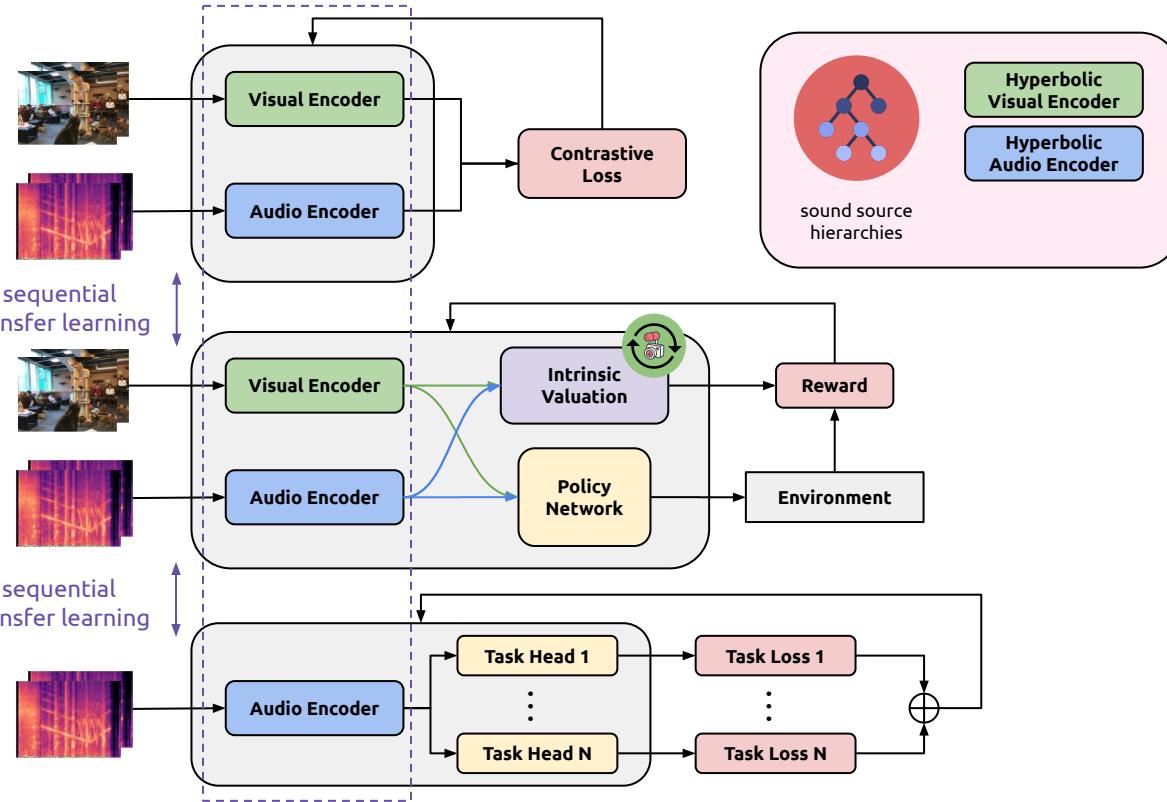
C. Gan *et al.*, "Noisy Agents: Self-supervised Exploration by Predicting Auditory Events," 2020.

V. Dean, S. Tulsiani, and A. Gupta, "See, Hear, Explore: Curiosity via Audio-Visual Association," *NeurIPS*, 2020.

# Putting it all together



**egocentric videos**



**embodied agents**



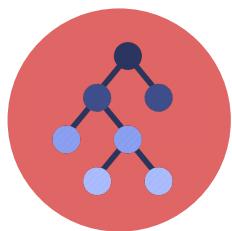
**multi-task  
datasets**

multi-task  
output structure

# In summary:

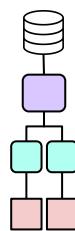
Incorporating natural structure is a promising path towards alleviating data-scarcity and improving robustness in machine listening models!

## parallel transfer learning (a.k.a. multi-task learning)



sound source  
hierarchies

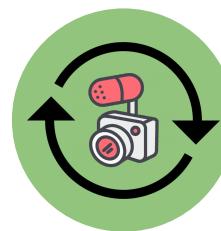
hierarchical structure  
improves model  
robustness!



multi-task output  
structure

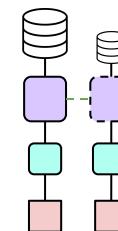
identifying and capitalizing  
on multi-task structure  
improves cross-task  
transfer!

## sequential transfer learning



multi-modal  
self-supervision

natural multi-modal  
correspondence provide  
structure without  
human annotations!



embodied agents

embodied navigation  
may provide structure  
better aligning with  
everyday experiences

# Thank you!



Juan Pablo Bello



Yao Wang



Ivan Selesnick



Michael Mandel



Dan Ellis



Ho-Hsiang Wu



Justin Salomon



Charlie Mydlarz



Mark Cartwright



Ana Elisa Méndez



Yu Wang



Vincent Lostanlen



Andrew Farnsworth



Benjamin Van Doren



Fatemeh Pishdadian



Magdalena Fuentes



Mahin Salman



Bea Steers



Gordon Wichern



Jonathan Le Roux



Rafael Valle



Bryan Catanzaro



Sivan Ding

+ many others!



**NYU**  
TANDON SCHOOL  
OF ENGINEERING



**CUNY** THE CITY  
UNIVERSITY OF NEW YORK

**Google**



The Cornell Lab  
of Ornithology

**NJIT**  
New Jersey Institute  
of Technology



Northwestern  
University



**MITSUBISHI**  
ELECTRIC

COLUMBIA  
UNIVERSITY



**L'ORÉAL** **BOSCH**



Olive