Robust Sound Event Detection in Acoustic Sensor Networks

Senior Research Scientist Music & Audio Research Lab / Center for Urban Science & Progress New York University

Work by MARL, SONYC & BirdVox





justinsalamon.com

Justin Salamon

@justinsalamon



justin.salamon@gmail.com



With support from:



NYU Music and Audio Research Lab (MARL)



justinsalamon.com

Acoustic Sensor Networks





justinsalamon.com

BIRDVOX











Anish Arora Co-PI



Carlos Bautista App Developer



Juan Pablo Bello Lead-PI



Mark Cartwright

Postdoc



Jason Cramer Masters Student



Harish Doraiswamy Res. Asst. Prof.



Graham Dove Postdoc



R. Luke DuBois Co-PI



Ben Esner

Senior Personnel

Peter Li PhD Student



Yitzchak Lockerman Postdoc

justinsalamon.com



Ana Elisa Mendez PhD Student



Fabio Miranda PhD Student



Charlie Mydlarz Senior Research Scientist



Oded Nov Co-PI





Yurii Piadyk PhD Student



Dhrubojyoti Roy PhD Student



Justin Salamon Senior Research Scientist



Ayanna Seals PhD Student



Mohit Sharma Asst. Research Scientist





Claudio Silva Co-PI



Yu Wang PhD Student



Ho-Hsiang Wu PhD Student

BirdVox team





Juan Pablo Bello NYU Andrew Farnsworth Cornell



Steve Kelling Cornell



Vincent Lostanlen Cornell / NYU



Kendra Oudyk McGill



Justin Salamon NYU



Why should we care about noise?





Washington Heights is among the city's noisiest neighborhoods, according to a recent study of 311 noise complaints. (Credit: Linda Rosier)

y NEW

Ø⁺ Ø £ ⊙

311 noise complaints on the rise in Washington Heights, Inwood

By Lisa L. Colangelo lisa.colangelo@amny.com January 30, 2018

New Yorkers are making a lot of noise about noise.

There were 1.6 million noise complaints made to 311 between 2010 and 2015, according to a new report released Monday by State Comptroller Thomas DiNapoli.

"Noise in New York City is a significant quality of life and public health concern," DiNapoli said. "The city ha noise code and should be commended for taking steps to better enforce local law, but there is more that ci



(/new-york-kids/sealand C imuni

includ Washington Heights an ^{al}and C imunity Boards 4 and 5

NEWS CITY LIFE

o ADD COMMENT 🖤 LOVE IT 🗖 SAVE I

Yikes! NYC's noisiest neighborhoods are no place for exhausted parents

h f ♥ G•
By Danielle Valente
Posted: Wednesday January 31 2018, 12:31pm





"I've had two years of absolute violation of my right to peace and quiet," said Mr. McIntosh, a television producer who has lived on the Upper East Side for more than five decades. "I think it's against the Geneva Conventions to have this much noise."

n at se. is l. l





SLEEP LOSS HEARING LOSS PR



justinsalamon.cc

Estimated 9 of 10 adults in NYC exposed to HARMFUL levels of NOISE



Over 3.4 MILLION complaints since 2003 [based on 311 data]



PRODUCTIVITY



STRESS



HOW DOES THE CITY CURRENTLY TACKLE NOISE?



justinsalamon.com



HOW DOES THE CITY CURRENTLY TACKLE NOISE?



justinsalamon.com



HOW DOES THE CITY CURRENTLY TACKLE NOISE?













justinsalamon.com



































Why do we care about classifying bird vocalizations?

White-throated_Sparrow : January 4



The migration monitoring puzzle





Previously on SONYC...















2 Years Later...

- BirdVox: over 6,000 hours of audio
- SONYC: over 26 years of audio
- New data, new challenges:
 - How do we label the data?
 - How can we best leverage the labeled data?
 - How can we leverage the unlabeled data?
- How do we make our models robust to different sensor locations?





justinsalamon.com

How do we label the data?



How do we label the data?

Scaper: A Library for Soundscape Synthesis and Augmentation J. Salamon, D. MacConnell, M. Cartwright, P. Li, and J. P. Bello. IEEE Workshop on Applications of Sig. Proc. to Audio and Acoustics (WASPAA), New Paltz, NY, USA, Oct. 2017.

Seeing sound: Investigating the effects of visualizations and complexity on crowdsourced audio annotations M. Cartwright, A. Seals, J. Salamon, A. Williams, S. Mikloska, D. MacConnell, E. Law, J. Bello, and O. Nov. Proceedings of the ACM on Human-Computer Interaction, 1(2), 2017.

Investigating the Effect of Sound-Event Loudness on Crowdsourced Audio Annotations M. Cartwright, J. Salamon, A. Seals, O. Nov, and J. P. Bello In IEEE Int. Conference on Acoustics, Speech, and Signal Processing (ICASSP), Calgary, Canada, Apr. 2018.

justinsalamon.com

Sound Event Detection (SED)





"Strong" sound event labels

- Strong labels:
 - Label (sound source)
 - Start time
 - End time





Crowdsourcing





The Audio Annotator

github.com/CrowdCurio/audio-annotator

Configured with the spectrogram visualization:





The Audio Annotator

Configured with the waveform visualization:





The Audio Annotator

Configured without a visualization:





Experiment

• 3 x 3 x 2 between-subjects factorial design:



<u>visualization</u>

• Soundscape examples: M0G0 M0G1

justinsalamon.com

max-polyphony

<u>gini-polyphony</u>





M2G0

M2G1



- Open source python library for soundscape synthesis
- Returns audio + annotation in JAMS format
- <u>github.com/justinsalamon/scaper</u> (pip install scaper)



justinsalamon.com

Figure 1: Block diagram of the Scaper synthesis pipeline.



Effect of Visualization on Quality and Speed of Annotations



Spectrogram \rightarrow higher-quality and faster annotations



Effect of Visualization on Task Learning



Expect even higher quality annotations after learning period



Effect of sound event loudness & scene complexity

Sound event SNR and overlap (polyphony) have direct impact on label recall, precision remains stable







Effect of Number of Annotators on Aggregate Annotation Quality



16 annotators captured 90% of gain in annotation quality, but 5 annotators is reasonable choice with respect to cost/quality trade-off



Annotating strong labels is slow...




How can we best leverage the labeled data? (can we get away without strong labels?)





How can we best leverage the labeled data? (can we get away without strong labels?)

Adaptive pooling operators for weakly labeled sound event detection B. Mcfee, J. Salamon, and J. P. Bello IEEE/ACM Transactions on Audio, Speech and Language Processing, 26(11): 2180–2193, Nov. 2018.





Sound Event Detection (SED)





Strong vs. weak labels

- Strong labels:
 - Label (sound source)
 - Start time
 - End time

- Weak labels (tags, e.g. AudioSet labels):
 - Label (sound source)
 - No timing information
 - Easier/faster/cheaper to curate!





Problem formulation

- Given dataset of tracks with **weak labels**:
 - [CLAPPING, SIREN]
- [CLAPPING, CAR HONK, SIREN, LAUGHTER]
 - [LAUGHTER]
- Train a SED model that outputs **strong labels**:



Sound Event Detection: If we had strong labels





Sound Event Detection: If we had strong labels



Reference (weak labels)



Multiple Instance Learning for Sound Event Detection





Training a MIL model for Sound Event Detection





Temporal pooling



- Brittle
- Gradient descent ignores most predictions
- Learning slow/sensitive to initialization

justinsalamon.com

- Lose specificity
- Only good events present most of the time



• Learning is well behaved



Temporal pooling



0.2*0.18 + 0.9*0.37 + 0.5*0.25 + 0.3*0.20 = 0.55

$$\hat{P}_s(Y \mid X) = \sum_{x \in X} \hat{p}(Y \mid x) \left(\frac{\exp \hat{p}(Y \mid x)}{\sum_{z \in X} \exp \hat{p}(Y \mid z)} \right)$$



- Weighted average
- Positive predictions get more weight
- Predicts like max, learns like mean



Softmax: bounded input, bounded output

- Inputs to softmax are probabilities in [0, 1]
- Outputs are positive weights, sum to 1
- Bounded input ⇔ Bounded output
- Weights converge to uniform for large bags







AutoPool

- Idea: remove input bound!
- Introduce temperature parameter: α
- Learn α jointly with model:

$$\hat{P}_{\alpha}(Y \mid X) = \sum_{x \in X} \hat{p}(Y$$

 $\cdot \left(\frac{\exp\left(\alpha \cdot \hat{p}(Y \mid x)\right)}{\sum_{z \in X} \exp\left(\alpha \cdot \hat{p}(Y \mid z)\right)} \right)$ (| x)

44

AutoPool

- $\alpha = 0 \Rightarrow$ mean pooling
- $\alpha = 1 \Rightarrow$ softmax pooling
- $\alpha \rightarrow \infty \Rightarrow \max pooling$
- $\alpha \rightarrow -\infty \Rightarrow$ min pooling
- Multi-label? each class gets its own α
- Model adapts to temporal extent of each class

$$\hat{P}_{\alpha}(Y \mid X) = \sum_{x \in X} \hat{p}(Y \mid x) \left(\frac{\exp\left(\alpha \cdot \hat{p}(Y \mid x)\right)}{\sum_{z \in X} \exp\left(\alpha \cdot \hat{p}(Y \mid z)\right)} \right)$$

• AutoPool variants: Constrained AutoPool (CAP) & Regularized AutoPool (RAP)



Experimental design

- Keep most of the model architecture **fixed**:
 - Deep convolutional neural network (CNN)
 - Based on audio subnetwork of Look, Listen and Learn (L3) architecture [Arandjelovic & Zisserman'17]
 - Output frame rate: ~3 frames/sec





Experimental design

- Compare **pooling** layers:
 - Max, mean, softmax, AutoPool, CAP, RAP ($\lambda = 10^{-2}$, 10⁻³, 10⁻⁴)
- Compare against strong training



justinsalamon.com



• Evaluate strong prediction (instance-level) accuracy: **F1**, P, R (0 = worst, 1 = best)



Experimental design: Datasets

- URBAN-SED [Salamon et al.'17]

 - 10 sound classes, mostly short events (~3 s)
- DCASE 2017 (Task 4) [Mesaros et al.'17]
 - ~50k subset of AudioSet [Gemmeke et al.'17]
 - 17 event classes, varying event durations (0-10 s)
- MedleyDB [Bittner et al.'15]
 - 122 songs -> 531 with remixing
 - 8 instrument classes, mostly long events (> 10s)









Results: URBAN-SED

	Strong (time-varying)				
Model	F_1	P	R	E_{\downarrow}	
Max	0.463	0.774	0.330	0.695	
Mean	0.408	0.280	0.751	2.10	
Soft-max	0.492	0.397	0.646	1.22	
RAP 10^{-2}	0.419	0.296	0.717	1.88	
RAP 10^{-3}	0.529	0.584	0.484	0.731	
RAP 10^{-4}	0.526	0.650	0.442	0.681	
CAP	0.533	0.622	0.466	0.696	
Auto	0.504	0.738	0.382	0.665	
Strong	0.551	0.693	0.458	0.642	



Results: DCASE 2017





Results: MedleyDB

Instance-level prediction (frames)

	Strong (time-varying)				
Model	F_1	P	R	E_{\downarrow}	
Max	0.437	0.875	0.292	0.719	
Mean	0.655	0.594			
Soft-max	0.662	0.668	0.658	0.524	
RAP 10^{-2}	0.659	0.604			
RAP 10^{-3}	0.673	0.638			
RAP 10^{-4}	0.622		0.530		
CAP	0.609		0.498		
Auto	0.528		0.386	0.636	
Strong	0.675	0.640	0.716	0.540	





Results: example

16384

4096

ΗZ 2048

1024

512

air conditioner car horn children playing dog bark drilling engine idling gun shot iackhammer jackhammer siren street music

air conditioner car horn children playing Estimate dog bark drilling engine idling gun shot jackhammer siren street music

Input

Reference

Model estimate



How can we leverage the unlabeled data?





How can we leverage the unlabeled data?

TBD

J. Cramer, H-H. Wang, J. Salamon, J. P. Bello Coming soon... 2019.





- The idea:

Large dataset of labeled audio





- Embedding
- Step 2: use embedding to train "downstream" model on target task
- Small dataset of labeled audio Deep embedding model



justinsalamon.com

Step 1: train embedding model on surrogate task

Prediction on surrogate task Deep embedding model Evaluate Update model parameters





- SoundNet [Aytar et al.'16]
 - Visual network pre-trained on large image datasets



justinsalamon.com

• Audio network trained to mimic output of visual network for Flickr videos



- VGGish [Hershey et al.'17]



justinsalamon.com

• Single audio network trained to predict labels on YouTube-8M dataset



- Look, Listen, and Learn (L3) [Arandjelovic & Zisserman'17]
 - Train model on the task of Audio-Visual Correspondence (AVC)
 - No labels required!





- Look, Listen, and Learn (L3) [Arandjelovic & Zisserman'17]
 - Train model on the task of Audio-Visual Correspondence (AVC)
 - No labels required!





Can we train our own (maybe better) L3 embedding?

- Questions:
 - Does an **audio-informed input representation** give a better embedding?
 - Is it important to use **matched audio domains** between embedding and target?
 - How much training data is enough training data for the embedding?
 - Does data augmentation still improve performance on target?
 - How does the resulting embedding compare to state-of-the-art embeddings?
 - SoundNet [Aytar et al.'16]
 - VGGish [Hershey et al.'17]





Experimental design (abridged)

- Step 1: train several variants of L3 embedding
 - Vary input representation, training data (content & amount)
- Step 2: evaluate embeddings on target task:
 - Multi-class sound classification using 2-layer MLP
 - Datasets:
 - UrbanSound8K (8732 clips, 10 classes) [Salamon et. al'14]
 - ESC-50 (**2000** clips, 40 classes) [Piczak'15]
 - DCASE 2013 (**200** clips, 10 classes) [Stowell et al.'15]
 - Compare to SoundNet and VGGish



L3: input representation





L3: matched vs mismatched training content





L3: embedding training data vs target performance



UrbanSound8K



Results

• Mel-based L3 vs SoundNet and VGGish:





Try these embeddings out... soon!

github.com/marl/openl3



How do we make our models robust to different sensor locations?
How do we make our models robust to different sensor locations? ...and to changing conditions in each location?



How do we make our models robust to different sensor locations? ...and to changing conditions in each location?

Per-Channel Energy Normalization: Why and how V. Lostanlen, J. Salamon, M. Cartwright, B. McFee, A. Farnsworth, S. Kelling, and J. P. Bello. IEEE Signal Processing Letters, in press, 2018.

Robust Sound Event Detection in Bioacoustic Sensor Networks V. Lostanlen, J. Salamon, A. Farnsworth, S. Kelling, and J. P. Bello. Coming soon... 2019













Challenge #2: nonuniform noise

Input representation: log-mel spectrogram?





Per-Channel Energy Normalization (PCEN)

pp. 5670–5674, New Orleans, LA, USA, Mar. 2017.

$$\mathbf{PCEN}(t,f) = \left(\frac{\mathbf{E}(t,f)}{\left(\varepsilon + \left(\mathbf{E}^{\text{time}} \ast \boldsymbol{\phi}\right)(t,f)\right)^{\alpha}} + \delta\right)^{r} - \delta^{r}$$

Adaptive gain control

justinsalamon.com

• Y. Wang, P. Getreuer, T. Hughes, R. F. Lyon, and R. A. Saurous. **Trainable** frontend for robust and far-field keyword spotting. In IEEE ICASSP,

Dynamic range compression



Per-Channel Energy Normalization (PCEN)

Signal Processing Letters, in press, 2018.



(b) Per-channel energy normalization (PCEN).

	indoor	outdoor
E	10 ⁻⁶	10 ⁻⁶
α	0.98	0.80
δ	10	2
r	0.5	0.25

justinsalamon.com

• V. Lostanlen, J. Salamon, M. Cartwright, B. McFee, A. Farnsworth, S. Kelling, and J. P. Bello. Per-Channel Energy Normalization: Why and how. IEEE



(a) Logarithmic transformation.







(b) Per-channel energy normalization (PCEN).



Per-Channel Energy Normalization (PCEN)





justinsalamon.com

(a) Logarithmic transformation.



Context-Adaptive Neural Networks



Representation of background (???)





Background (auxiliary) Features

Summarize background at scale T = 30 minutes: energy quantiles (median, quartiles, centiles, etc.)

Auxiliary features describe noise only, not the class of interest.

mel frequency





Context-Adaptive Neural Networks

static sigmoid layer:

dynamic filter network:

mixture of experts:

adaptive threshold:

$$y(t) = \sigma \left(\sum_{k} \langle w_{k} | x_{k}(t) \rangle + b \right)$$
$$y(t) = \sigma \left(\sum_{k} \langle w_{k}(t) | x_{k}(t) \rangle + b \right)$$
$$y(t) = \sigma \left(\sum_{k} \alpha_{k}(t) \langle w_{k} | x_{k}(t) \rangle + b \right)$$
$$y(t) = \sigma \left(\sum_{k} \langle w_{k} | x_{k}(t) \rangle + b(t) \right)$$



Context-Adaptive Neural Networks





- Spectral flux: 15%
- CNN: 56%
- CNN + aug: 62%
- CNN + PCEN + aug: 66%
- CNN + PCEN + CA + aug: 72%



Robust Sound Event Detection in Acoustic Sensor Networks

- How do we label the data?
 - Crowdsourcing: <u>github.com/CrowdCurio/audio-annotator</u>
 - Synthesis: <u>github.com/justinsalamon/scaper</u>
- How can we best leverage the labeled data?
 - Multiple instance learning + AutoPool: <u>github.com/marl/autopool</u>
- How can we leverage the unlabeled data?
 - Deep audio embeddings: <u>github.com/marl/openl3</u>
- How do we make our models robust to different & changing sensor locations?
 - PCEN: <u>github.com/librosa/librosa</u>

justinsalamon.com

Context-adaptive networks: <u>github.com/BirdVox/bv_context_adaptation</u>

Thanks!

