Learning spectro-temporal representation of complex sounds with parameterized neural networks

Rachid Riad, Julien Karadayi, Anne-Catherine Bachoud-Levi and Emmanuel Dupoux

https://arxiv.org/abs/2103.07125

















# Examples of auditory inputs



●

# Examples of auditory inputs

- 1. People chatting inside
- 2. Car alarm
- 3. Bird

# Examples of auditory inputs

- 1. People chatting inside
- 2. Car alarm
- 3. Bird





How do we obtain meaningful information from sounds?



How do we obtain meaningful information from sounds?

Auditory Neuroscience

How sounds are represented in the brain?





How do we obtain meaningful information from sounds?

Machine Listening

Can we design an algorithm to deal with sounds?



How do we obtain meaningful information from sounds?



How do we obtain meaningful information from sounds?



Not so much in practice!!

The Ear (Human Anatomy)















Α









Arnal et al. 2015

Santoro et al. 2017



- Does not account for **behaviour**
- Does not correlate much with the brain

# Deep Learning models as models of audition

Human level performance reached by Deep Neural Networks based on the spectrogram



**Convolution layers** 

0	1	1	1	0	.0	0	*********									
0	0	1	1,0	1	Q	0	**********		*****			1	4	3	4	1
0	0	0	1,	1	1	0		1	0	1		1	2	4	3	3
0	0	0	1	1.	.0	0	******	0	1	0	and the state	1	2	3	4	1
0	0	1	1	0	0	0	*********	1	0	1	in and in a second second	1	3	3	1	1
0	1	1	0	0	0	0			9 K			3	3	1	1	0
1	1	0	0	0	0	0						2	2		2 1	5

**Recurrent layers** 



 $z_t = \sigma \left( W_z \cdot [h_{t-1}, x_t] \right)$   $r_t = \sigma \left( W_r \cdot [h_{t-1}, x_t] \right)$   $\tilde{h}_t = \tanh \left( W \cdot [r_t * h_{t-1}, x_t] \right)$   $h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$ 18

## Deep Learning models as models of audition

Human level performance reached by Deep Neural Networks based on the spectrogram



for neuroscientists and psychologists

**Convolution layers** 

0	1	1	1	0	.0	0										
0	0	1	1,0	1	Q,	0	**********					1	4	3	4	1
0	0	0	1,	1,0	1,	0		1	0	1		1	2	4	3	3
0	0	0	1	·1.	.0	Ö	******	0	1	0		1	2	3	4	1
0	0	1	1	0	0	0	********	1	0	1	in a second de la seconde d	1	3	3	1	1
0	1	1	0	0	0	0					24 March 19	3	3	1	1	0
1	1	0	0	0	0	0						2 22	2 10	5 - 53	ā 14	5

#### **Recurrent layers**



 $z_t = \sigma \left( W_z \cdot [h_{t-1}, x_t] \right)$   $r_t = \sigma \left( W_r \cdot [h_{t-1}, x_t] \right)$   $\tilde{h}_t = \tanh \left( W \cdot [r_t * h_{t-1}, x_t] \right)$  $h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_{t_{19}}$ 

# Deep Learning models as models of audition

Human level performance reached by Deep Neural Networks based on the spectrogram

Not **interpretable** for neuroscientists and psychologists



**Convolution layers** 

0	1	1	1	•0	.0	0										
0	0	1	1,0	1,	Q,	Ö.	**********					1	4	3	4	1
0	0	0	1,	1	1	0		1	0	1		1	2	4	3	3
0	0	0	1	1.	.0	0	******	0	1	0		1	2	3	4	1
0	0	1	1	0	0	0	*********	1	0	1	in and in the second second	1	3	3	1	1
0	1	1	0	0	0	0					and the second second	3	3	1	1	0
1	1	0	0	0	0	0										5

#### **Recurrent layers**









#### Like the Gabor patterns found in Vision processing like in V1



#### Like the Gabor patterns found in Vision processing like in V1





$$igstarrow g_k(t,f) = s_k(t,f) w_k(t,f) \ s_k(t,f) = rac{1}{2\pi\sigma_t\sigma_f} e^{-rac{1}{2}\left(rac{t^2}{\sigma_t^2} + rac{f^2}{\sigma_f^2}
ight)} \ w_k(t,f) = e^{j(2\pi(\omega_k t + \Omega_k f))} \ w_k(t,f) = e^{j(2\pi(F_k R_{\gamma_k}))} \ R_{\gamma_k} = t\cos(\gamma_k) + f\sin(\gamma_k)$$



$$egin{aligned} g_k(t,f) &= s_k(t,f) w_k(t,f) \ s_k(t,f) &= rac{1}{2\pi\sigma_t\sigma_f} e^{-rac{1}{2}\left(rac{t^2}{\sigma_t^2} + rac{f^2}{\sigma_f^2}
ight)} \ w_k(t,f) &= e^{j(2\pi(\omega_k t + \Omega_k f))} \ w_k(t,f) &= e^{j(2\pi(F_k R_{\gamma_k}))} \ R_{\gamma_k} &= t\cos(\gamma_k) + f\sin(\gamma_k) \end{aligned}$$

**The idea**: Stay in the gabor Domain for the convolution, but find the relevant ones



**The idea**: Stay in the gabor Domain for the convolution, but find the relevant ones



27

How to implement the idea: Integrate convolutions into a neural network and backprop only a subset of parameters



$$egin{aligned} g_k(t,f) &= s_k(t,f) w_k(t,f) \ s_k(t,f) &= rac{1}{2\pi\sigma_t\sigma_f} e^{-rac{1}{2}\left(rac{t^2}{\sigma_t^2} + rac{f^2}{\sigma_f^2}
ight)} \ w_k(t,f) &= e^{j(2\pi(\omega_k t + \Omega_k f))} \ w_k(t,f) &= e^{j(2\pi(F_k R_{\gamma_k}))} \ R_{\gamma_k} &= t\cos(\gamma_k) + f\sin(\gamma_k) \end{aligned}$$

Learn only  $(\sigma_t, \sigma_f, R_{\gamma_k}, F_k)$ 

4 params instead of 111\*9

**How to implement the idea**: Integrate convolutions into a neural network and backprop only a subset of parameters



$$egin{aligned} g_k(t,f) &= s_k(t,f) w_k(t,f) \ s_k(t,f) &= rac{1}{2\pi\sigma_t\sigma_f} e^{-rac{1}{2}\left(rac{t^2}{\sigma_t^2} + rac{f^2}{\sigma_f^2}
ight)} \ w_k(t,f) &= e^{j(2\pi(\omega_k t + \Omega_k f))} \ w_k(t,f) &= e^{j(2\pi(F_k R_{\gamma_k}))} \ R_{\gamma_k} &= t\cos(\gamma_k) + f\sin(\gamma_k) \end{aligned}$$

L

<sub>earn only</sub> 
$$(\sigma_t, \sigma_f, R_{\gamma_k}, F_k)$$

4 params instead of 111\*9

# Our pipeline for audio processing



# Our pipeline for audio processing



# **Engineering results**

**Voice Activity Detection** 

**Speaker Identification** 

**Bird Call type classification** 

Audio scene classification

State-of-the-art on two datasets

Behind topline, above baseline

Behind topline, above baseline

Close to topline

# Qualitative analysis of filters



# Quantitative analysis of filters



#### Quantitative Analysis of spectro-temporal



#### Quantitative Analysis of spectro-temporal





#### Quantitative Analysis of spectro-temporal



# Quantitative analysis of filters



## What is next?

1. Harmonicity



### What is next?

- 1. Harmonicity
- 2. Tonotopic Organization for the Modulation



Figure 5. Organization of MTF in STG, LF, renears ductare centroid spectrated from all MTF across all participants. Their respective 59% contours are shown holes. The orerall lunning within an individual MTE centroid is farly well obtained in Star for MTF and the Centroid types span modulation space from high spectral/low temporal epicidual of the Star MTF and the method well are shown by the 59% contours. *B*, Gowg MTF may, The may respectents the average MTF across participants at each 37G position. Only lo catalox with  $\ge 2.017S$  contributing to the average are included. Lead MTF with the mesh is close coded by the distoret membership. (Catal Center MTF and Terror Star Center MTF and Terror Star Centroid (Terror Star Center MTF) and the method of the star center of the star Center MTF and the star center of the sta