

Renormaliser le spectrogramme : pourquoi ? comment ?

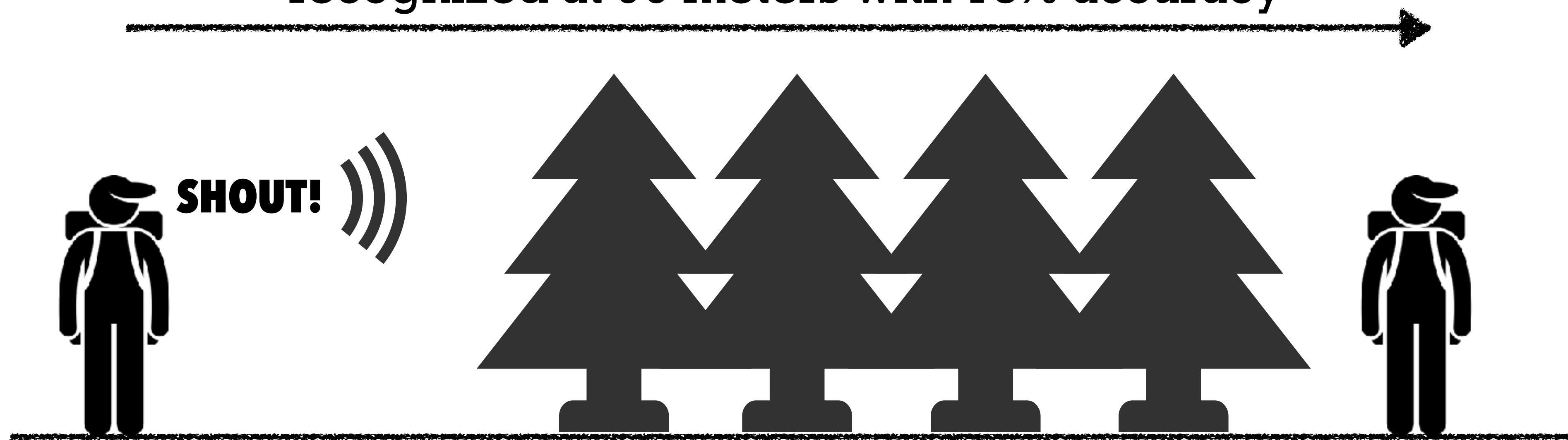


Vincent Lostanlen
LS2N, Centrale Nantes, CNRS

Hearing at a distance

[Meyer et al. 2018]

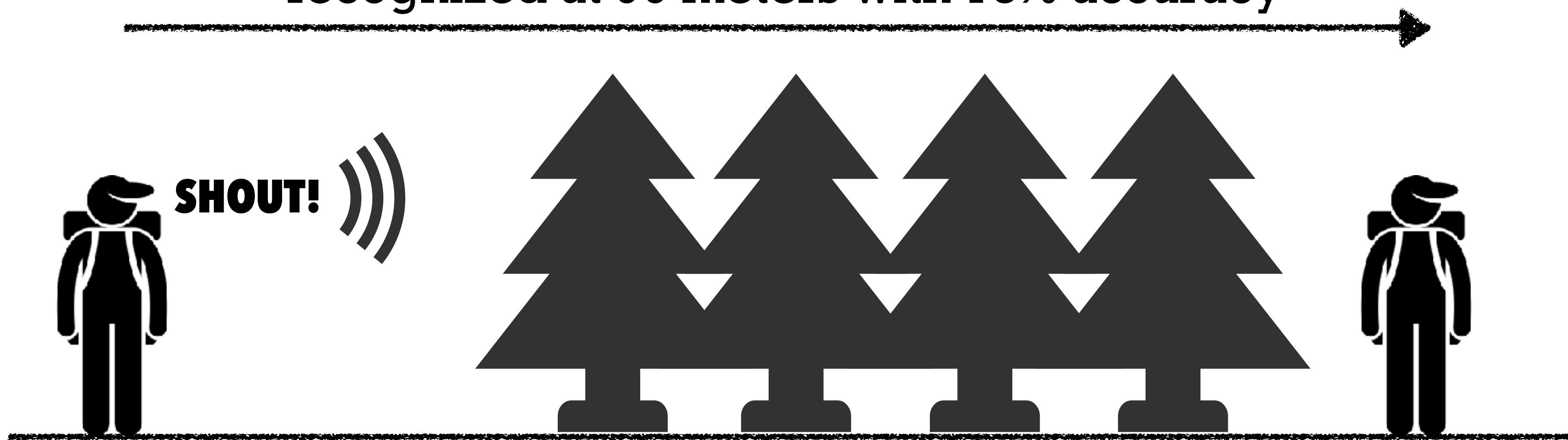
recognized at 90 meters with 75% accuracy



Hearing at a distance

[Meyer et al. 2018]

recognized at 90 meters with 75% accuracy

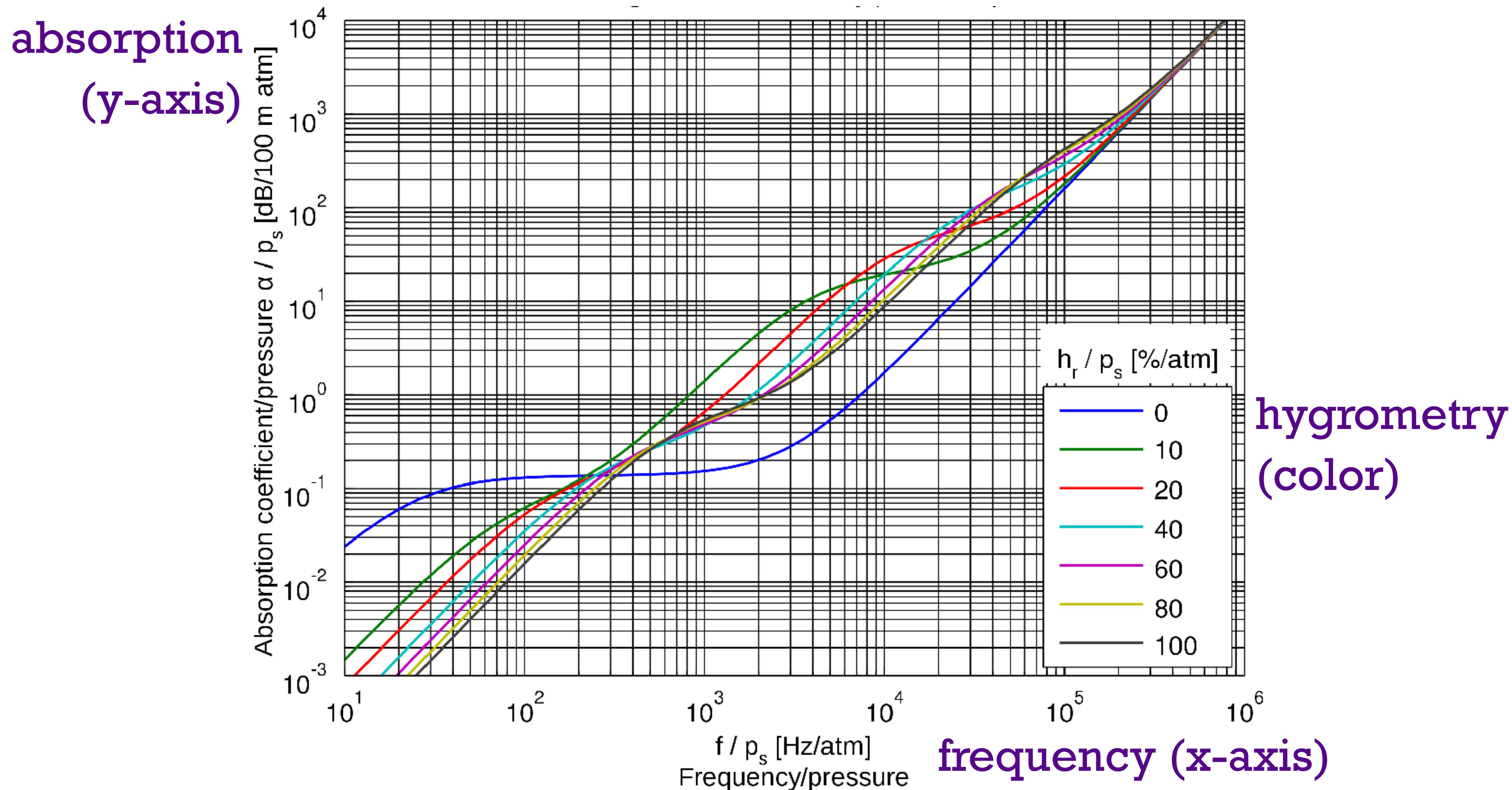


Yet, long-distance recognition is difficult for machines.

Why?

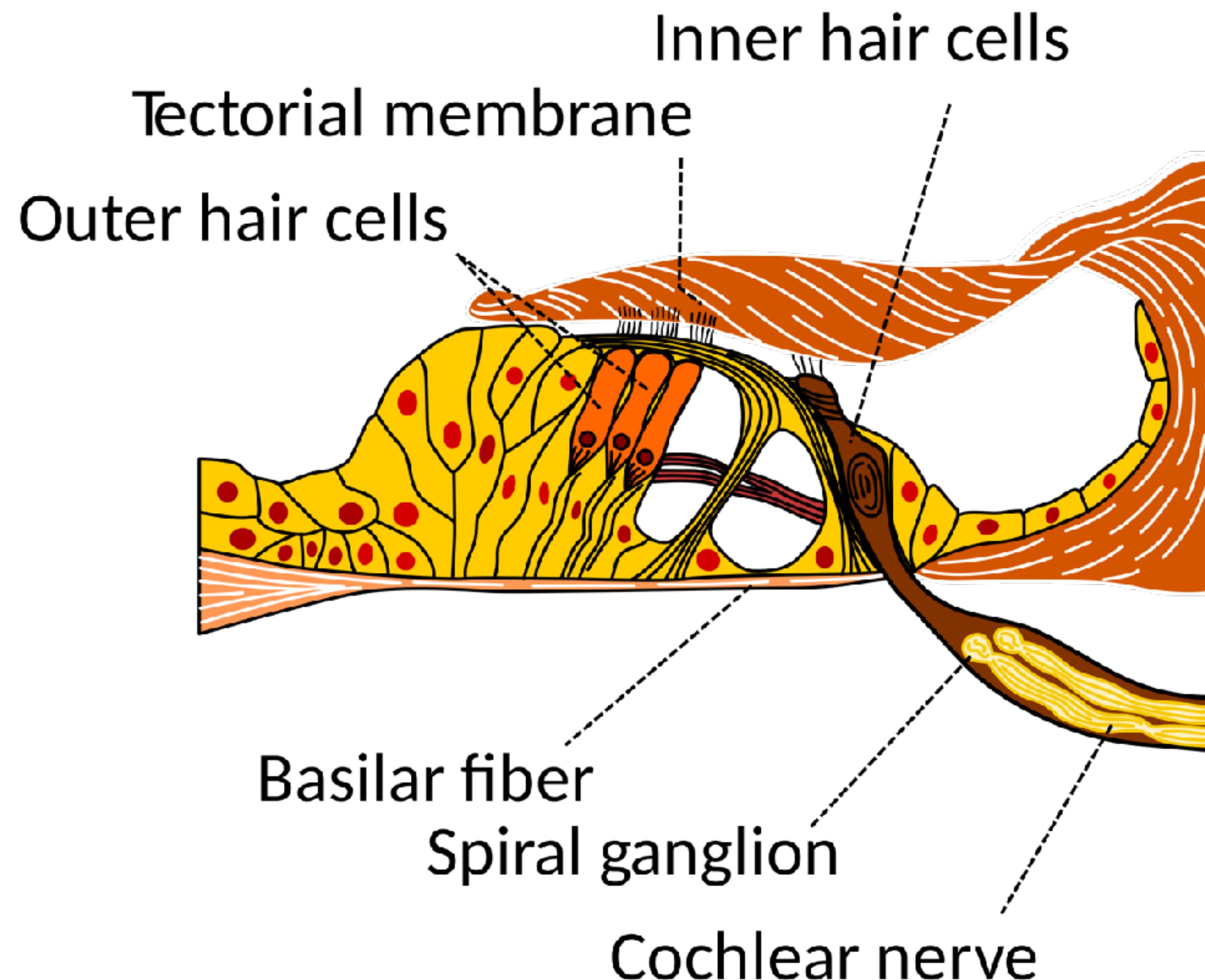
Atmospheric absorption

[Bass et al. 1984]



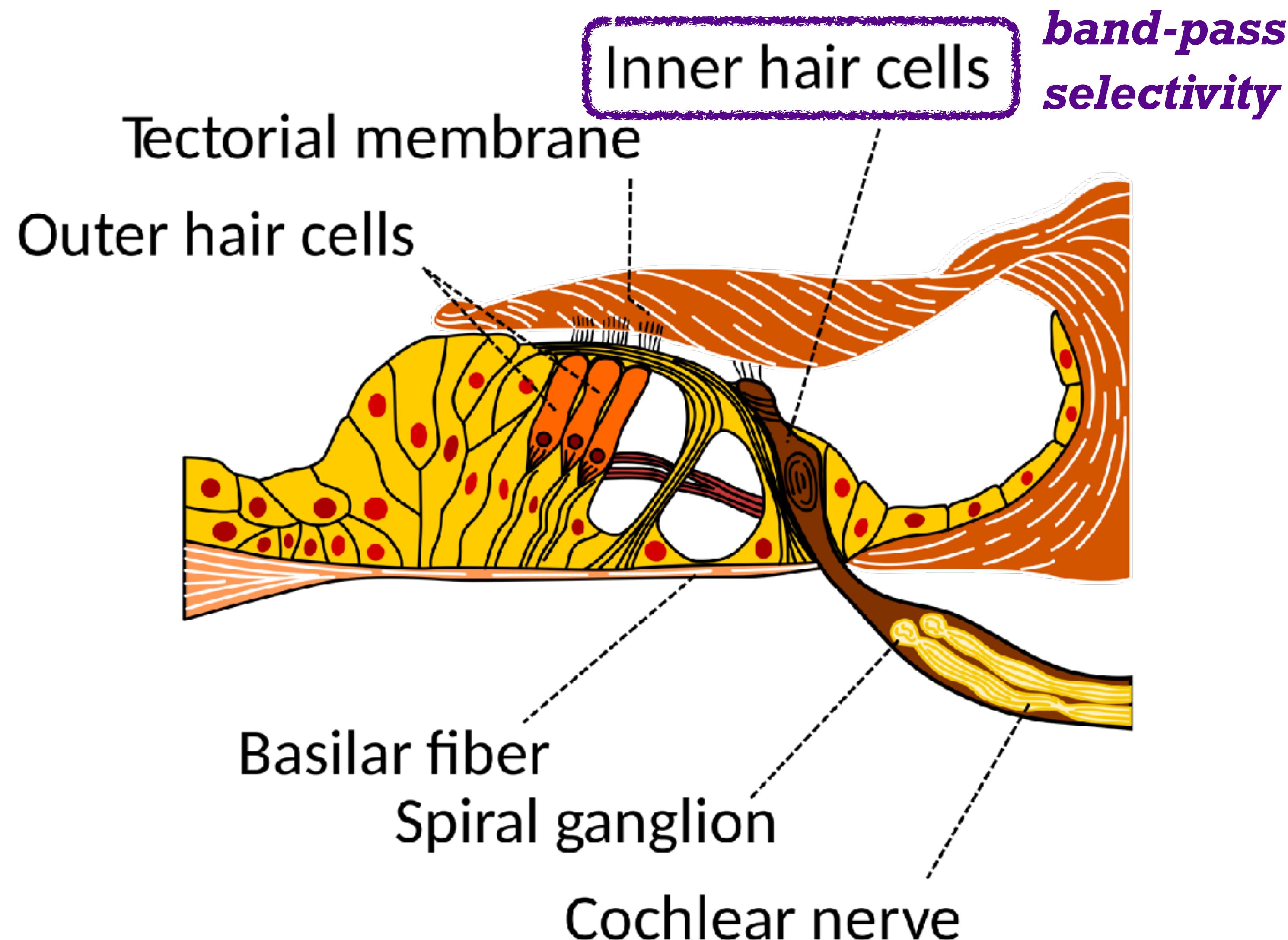
Auditory neurophysiology 101

[Mann and Kelley 2011, Robles et al. 2001]



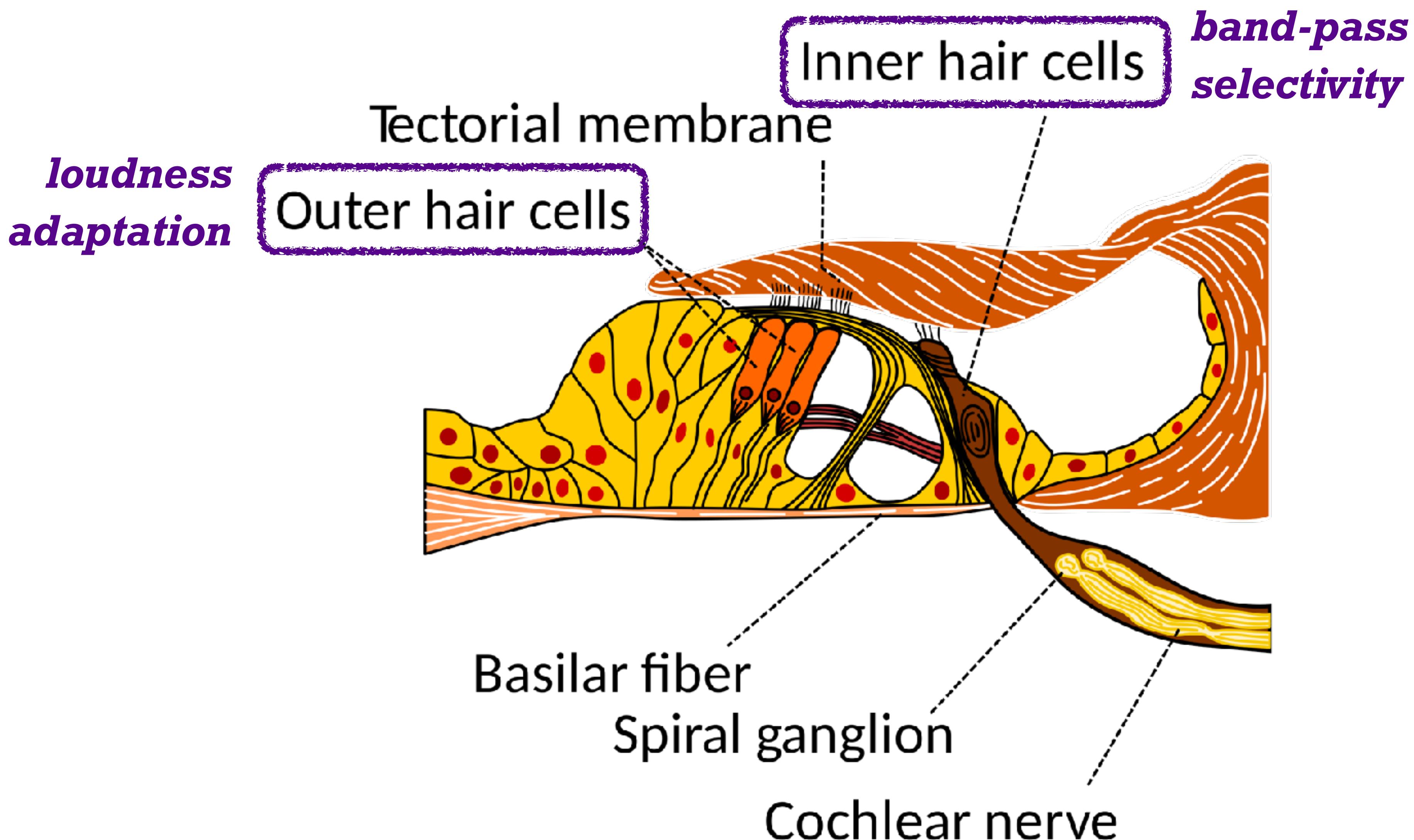
Auditory neurophysiology 101

[Mann and Kelley 2011, Robles et al. 2001]

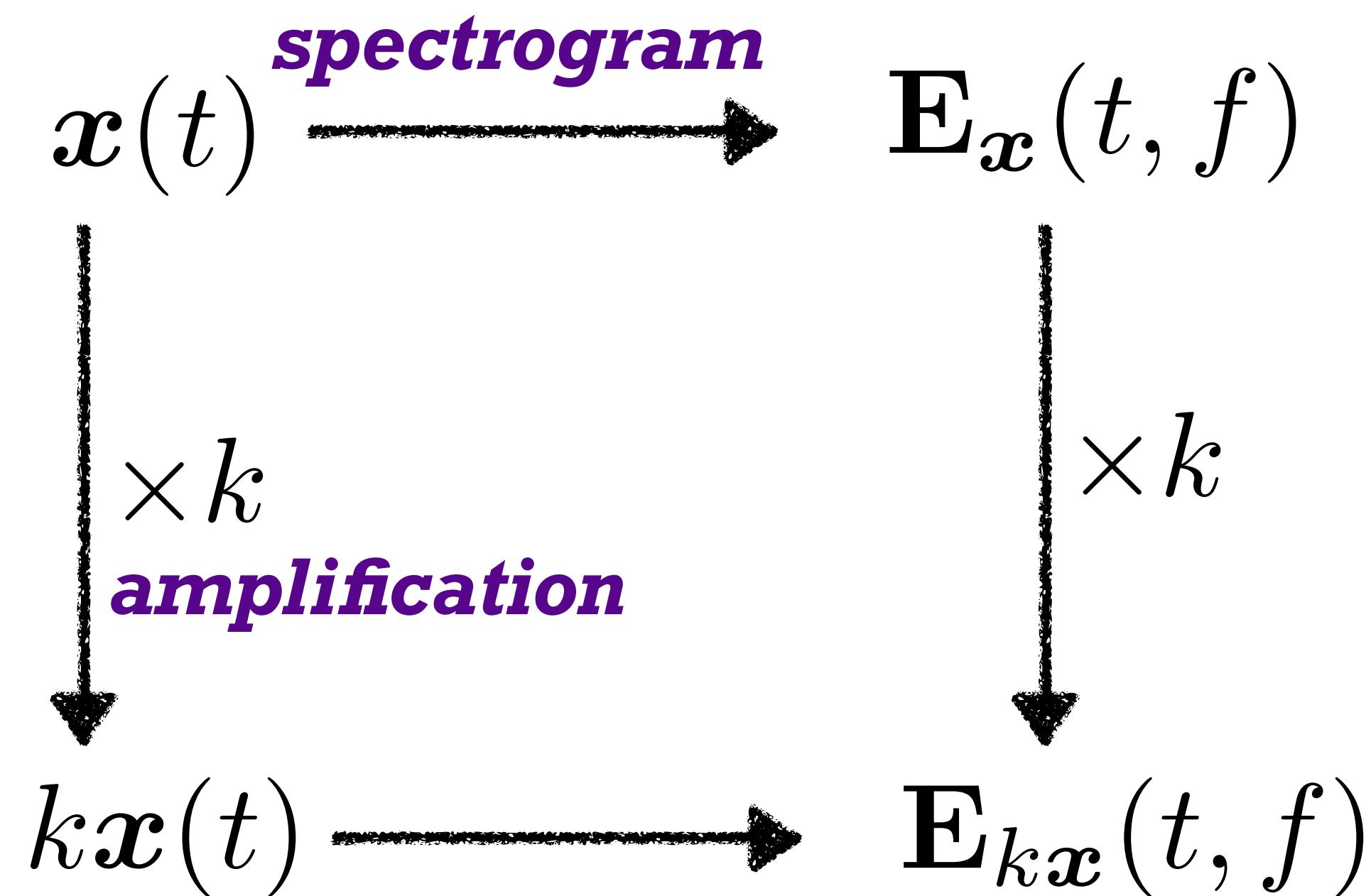


Auditory neurophysiology 101

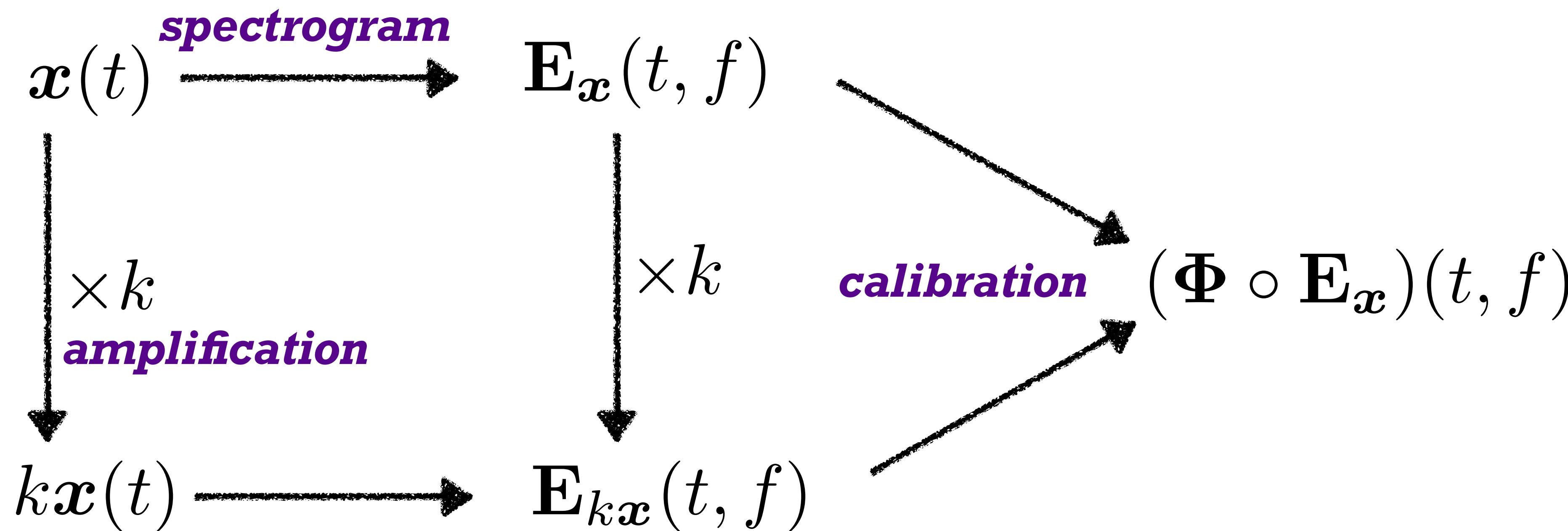
[Mann and Kelley 2011, Robles et al. 2001]



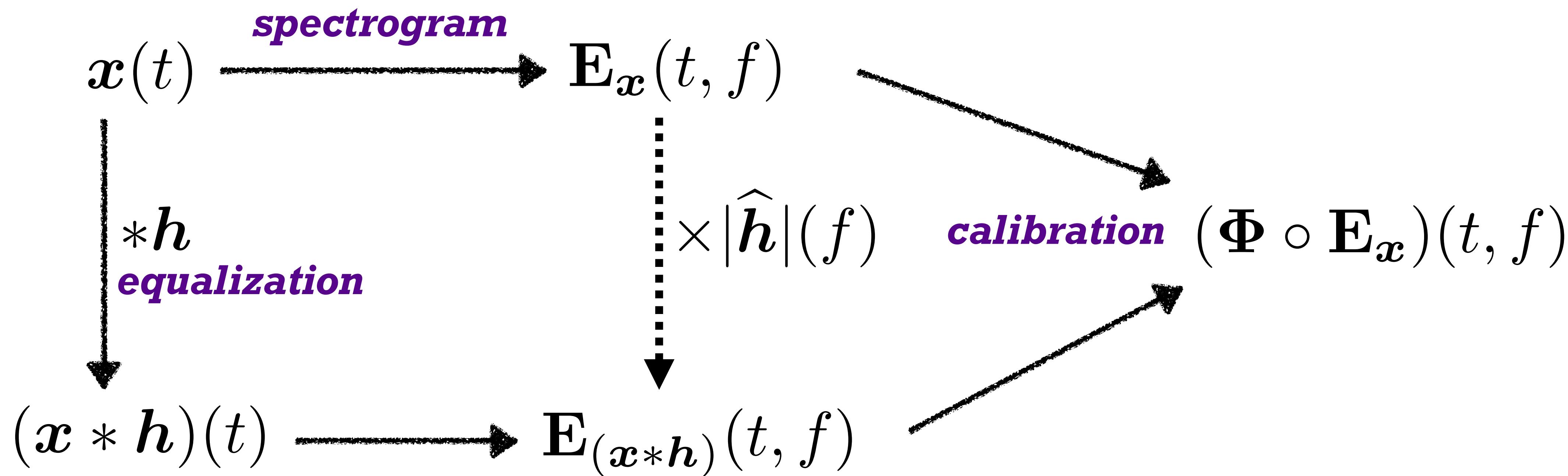
Invariance to amplitude



Invariance to amplitude



Invariance to equalization



Per-Channel Energy Normalization

[Wang et al. 2017, Lostanlen et al. 2019]

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

Per-Channel Energy Normalization

[Wang et al. 2017, Lostanlen et al. 2019]

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

*hyperparameters
(predefined or learned)*

Per-Channel Energy Normalization

[Wang et al. 2017, Lostanlen et al. 2019]

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

adaptive gain control

dynamic range compression

Per-Channel Energy Normalization

[Wang et al. 2017, Lostanlen et al. 2019]

$$\phi_T(t) = \exp(-t/T) \mathbf{1}(t > 0)$$

low-pass filter

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

Typically:

- **r** ~ 0.33
- **delta** ~ 0.5
- **epsilon** ~ 0
- **alpha** ~ 1
- **T** ~ 400 ms

*adaptive
gain control*

*dynamic range
compression*

Simplified PCEN

Under the simplest parametrization, PCEN boils down to

$$(\Phi \circ \mathbf{E}_x)(t, f) = \frac{\mathbf{E}_x(t, f)}{\mathbf{M}_x(t, f)}$$

where $\mathbf{M}_x(t, f)$ is a moving average of $\mathbf{E}_x(t, f)$.

Simplified PCEN

Under the simplest parametrization, PCEN boils down to

$$(\Phi \circ \mathbf{E}_x)(t, f) = \frac{\mathbf{E}_x(t, f)}{\mathbf{M}_x(t, f)}$$

where $\mathbf{M}_x(t, f)$ is a moving average of $\mathbf{E}_x(t, f)$.

- If $\mathbf{E}_x(t, f)$ is a “noise” region: $\mathbf{E}_x(t, f) \approx \mathbf{M}_x(t, f)$
and thus $(\Phi \circ \mathbf{E}_x)(t, f) \approx 1$ on average

Simplified PCEN

Under the simplest parametrization, PCEN boils down to

$$(\Phi \circ E_x)(t, f) = \frac{E_x(t, f)}{M_x(t, f)}$$

where $M_x(t, f)$ is a moving average of $E_x(t, f)$.

- If $E_x(t, f)$ is a “noise” region: $E_x(t, f) \approx M_x(t, f)$
and thus $(\Phi \circ E_x)(t, f) \approx 1$
- If $E_x(t, f)$ is a “transient” region: $E_x(t, f) \gg M_x(t, f)$
and thus $(\Phi \circ E_x)(t, f) \gg 1$

Simplified PCEN

Under the simplest parametrization, PCEN boils down to

$$(\Phi \circ E_x)(t, f) = \frac{E_x(t, f)}{M_x(t, f)}$$

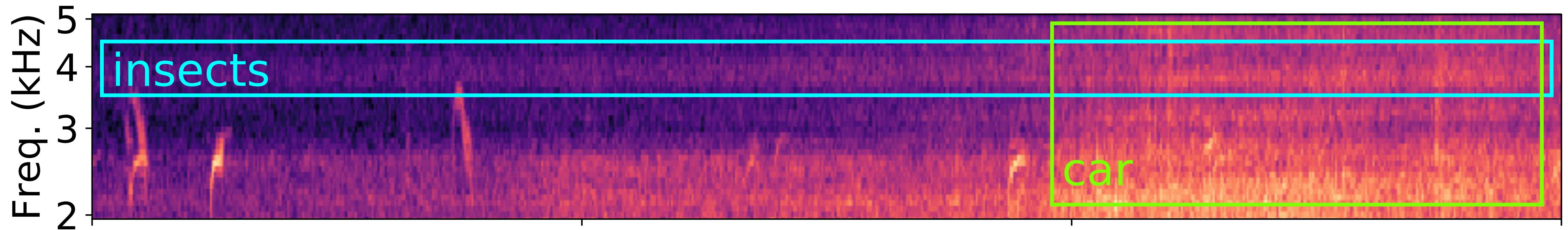
where $M_x(t, f)$ is a moving average of $E_x(t, f)$.

- If $E_x(t, f)$ is a “noise” region: $E_x(t, f) \approx M_x(t, f)$
and thus $(\Phi \circ E_x)(t, f) \approx 1$
- If $E_x(t, f)$ is a “transient” region: $E_x(t, f) \gg M_x(t, f)$
and thus $(\Phi \circ E_x)(t, f) \gg 1$

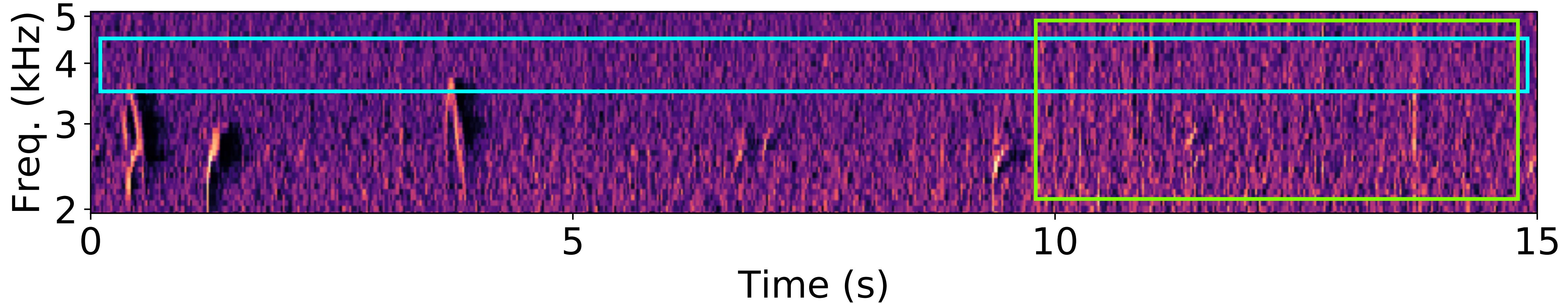
PCEN in practice

PCEN enhances the visual contrast between events and background.

before PCEN



after PCEN

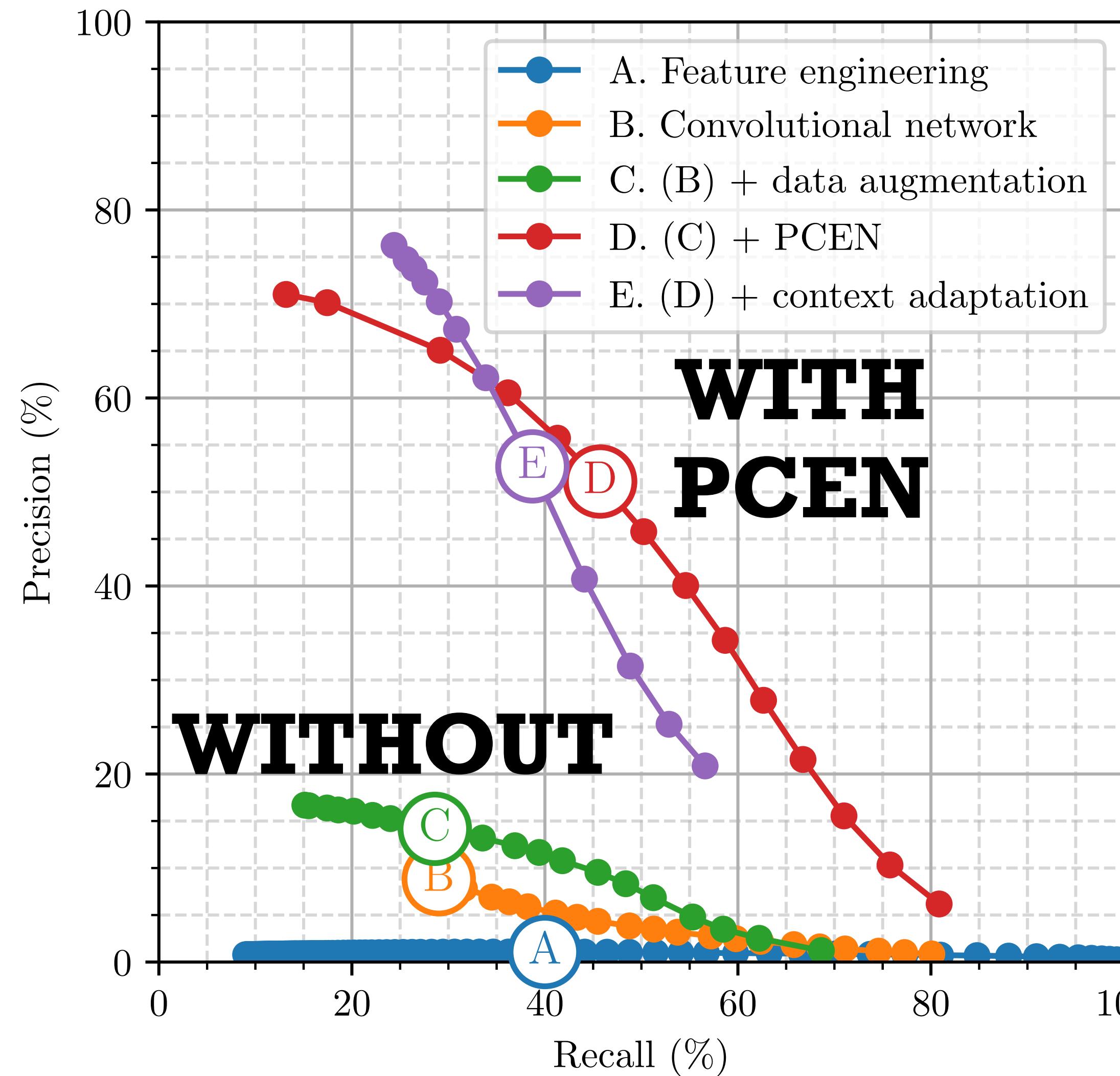


Data: BirdVox project

PCEN + deep learning

[Lostanlen et al. 2022]

Evaluated on 300h hours of annotated audio from a sensor network.



PCEN plays a crucial role
in statistical generalization.

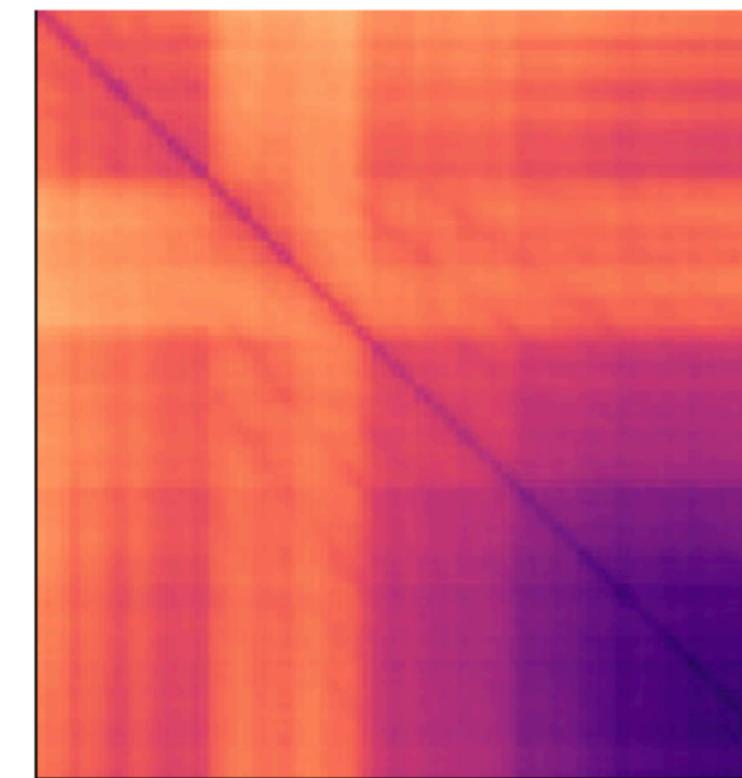
Models A, B, C: without PCEN
Models D, E: with PCEN

how to justify this?

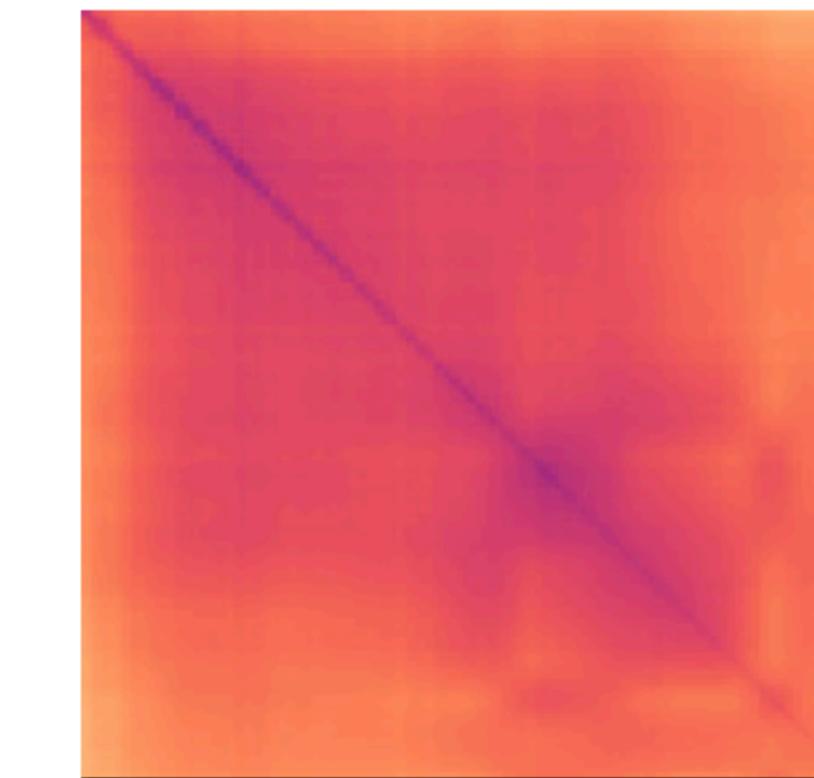
PCEN decorrelates subbands

[Lostanlen et al. 2019a]

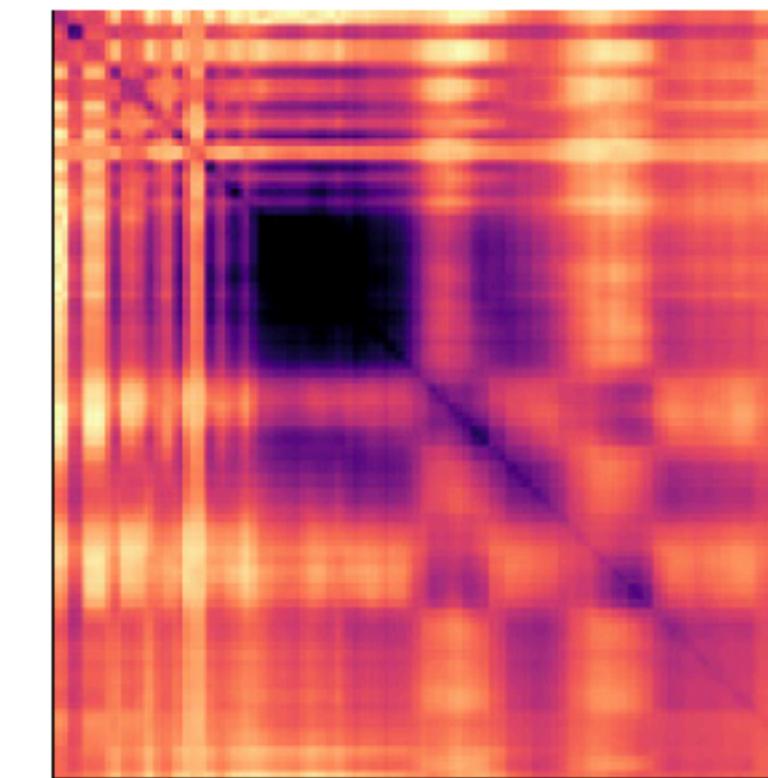
SONYC



DCASE 2013 SC

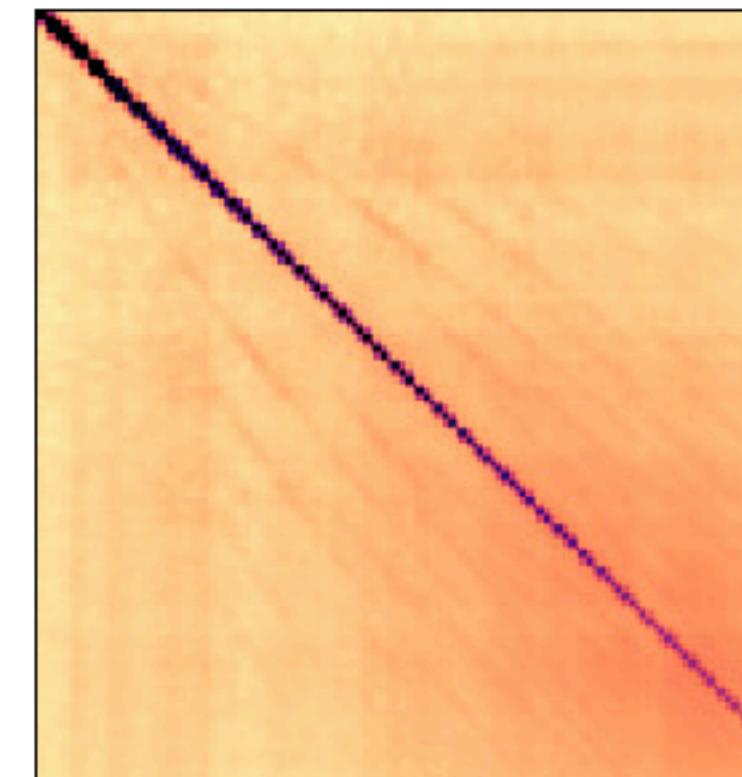


BirdVox

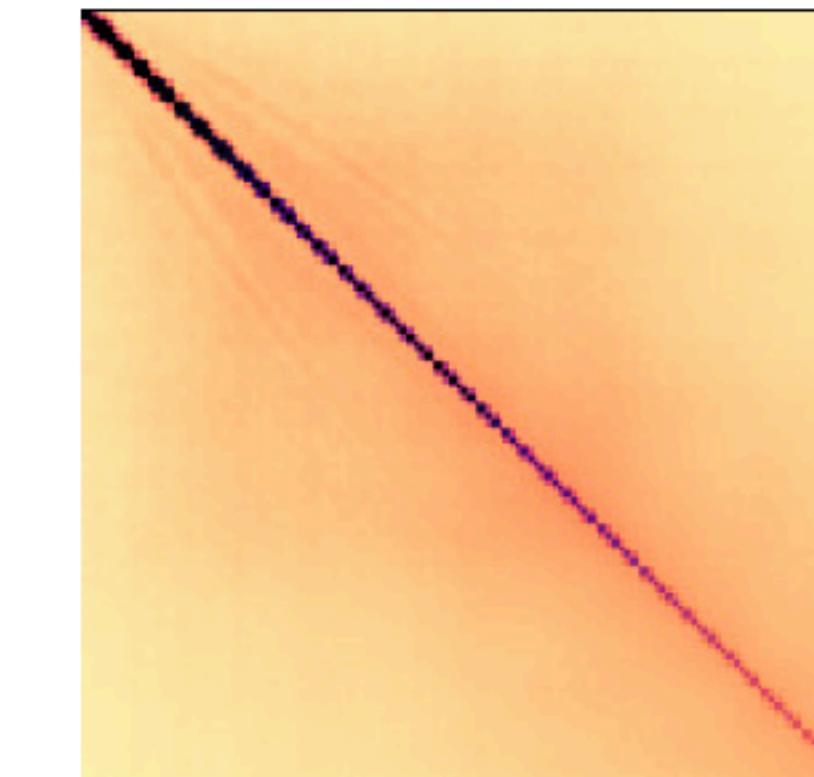


(a) Logarithmic transformation.

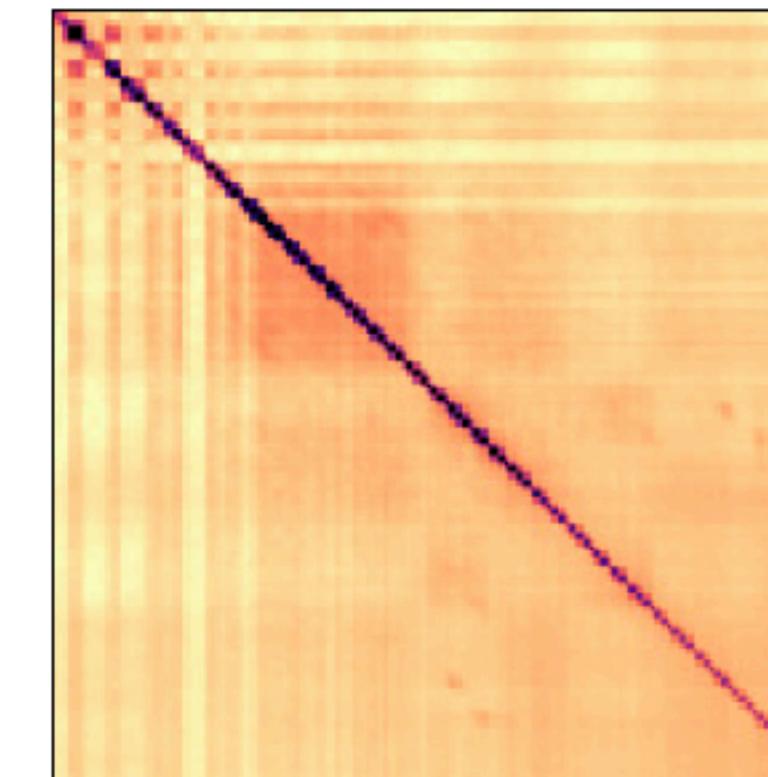
(urban noise)



(indoor noise)



(rural noise)



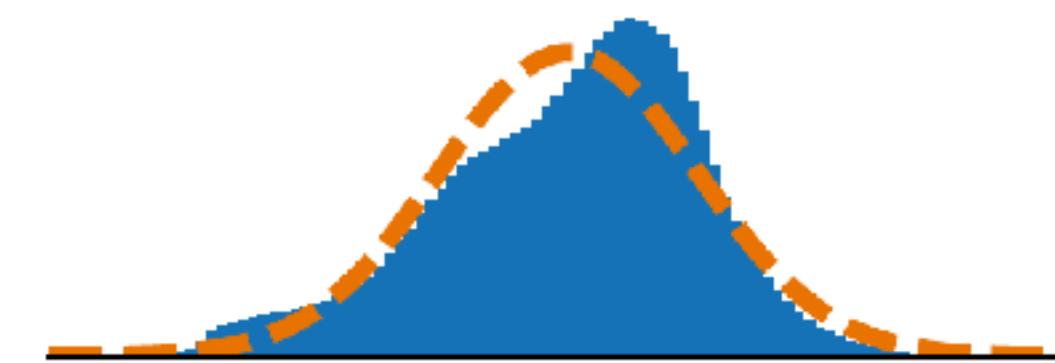
(b) Per-channel energy normalization (PCEN).

PCEN Gaussianizes magnitudes

[Lostanlen et al. 2019a]

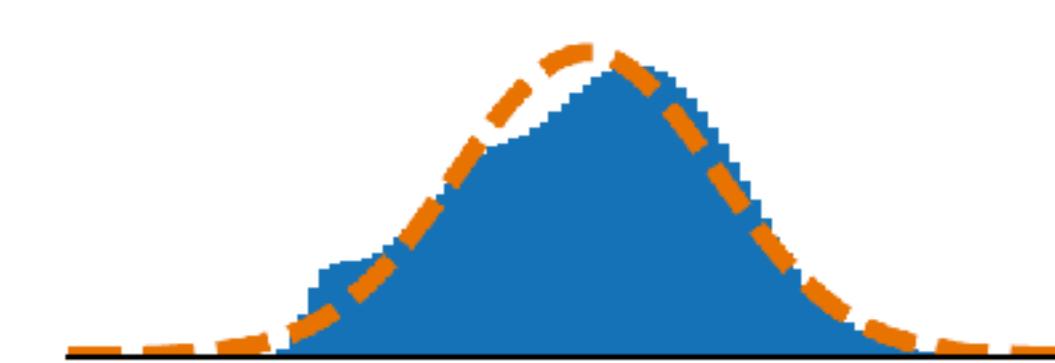
(urban noise)

SONYC



(indoor noise)

DCASE 2013 SC

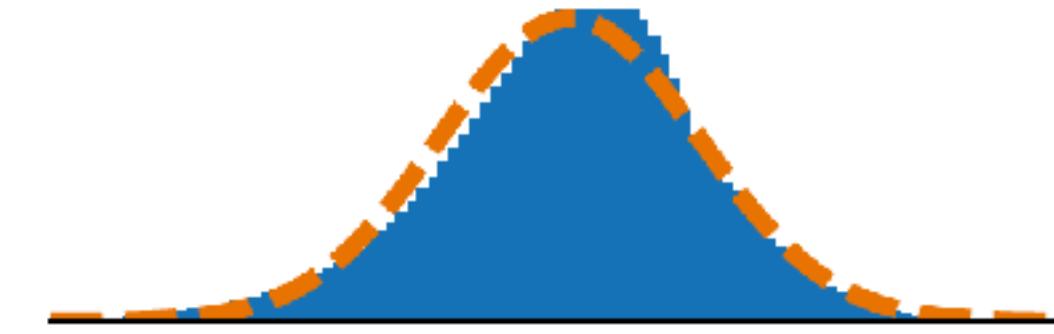
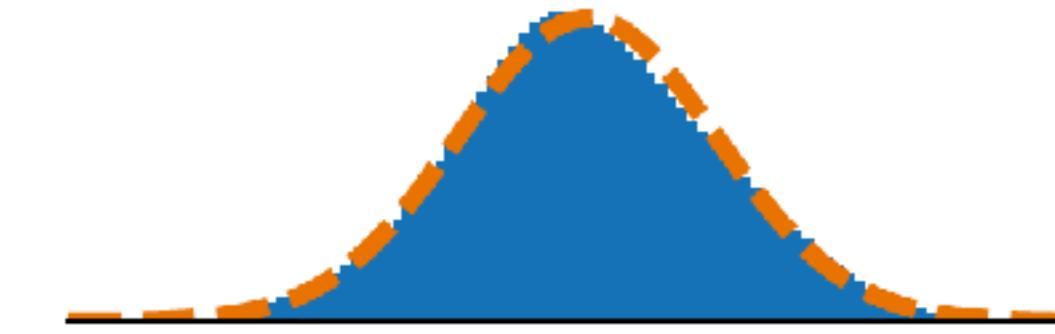
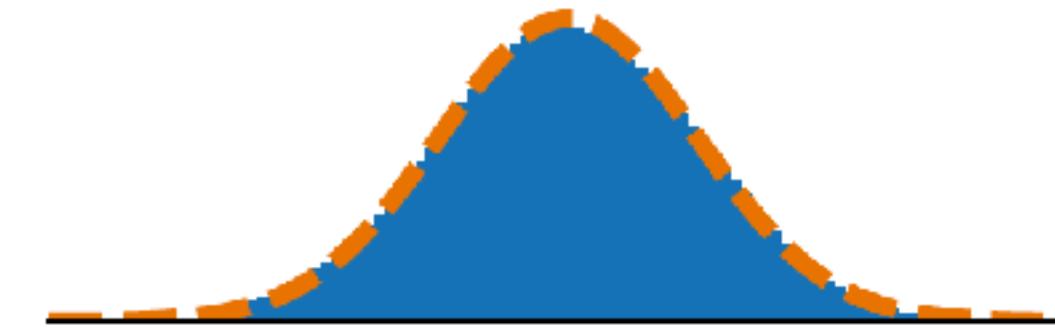


(rural noise)

BirdVox



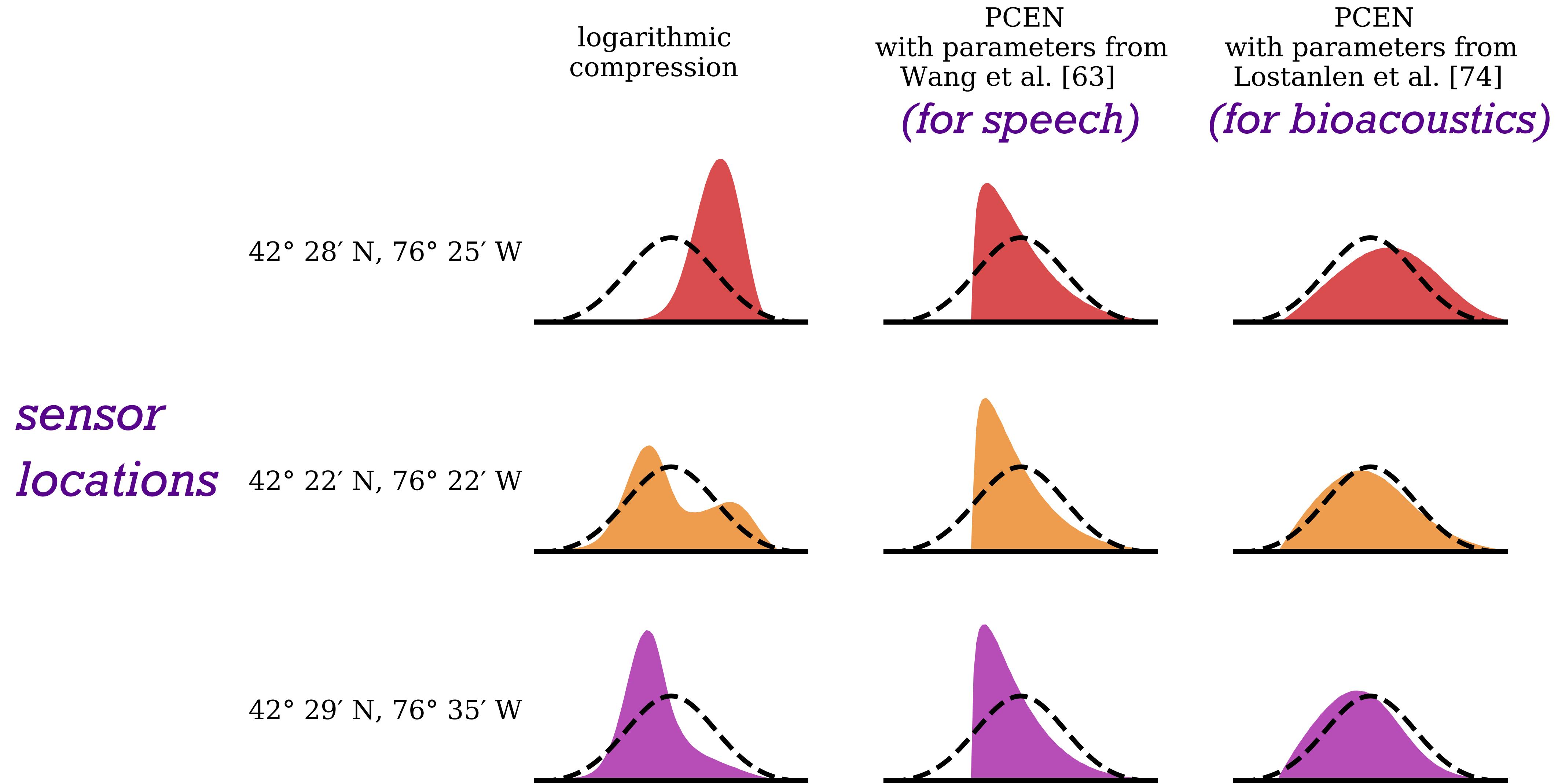
(a) Logarithmic transformation.



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

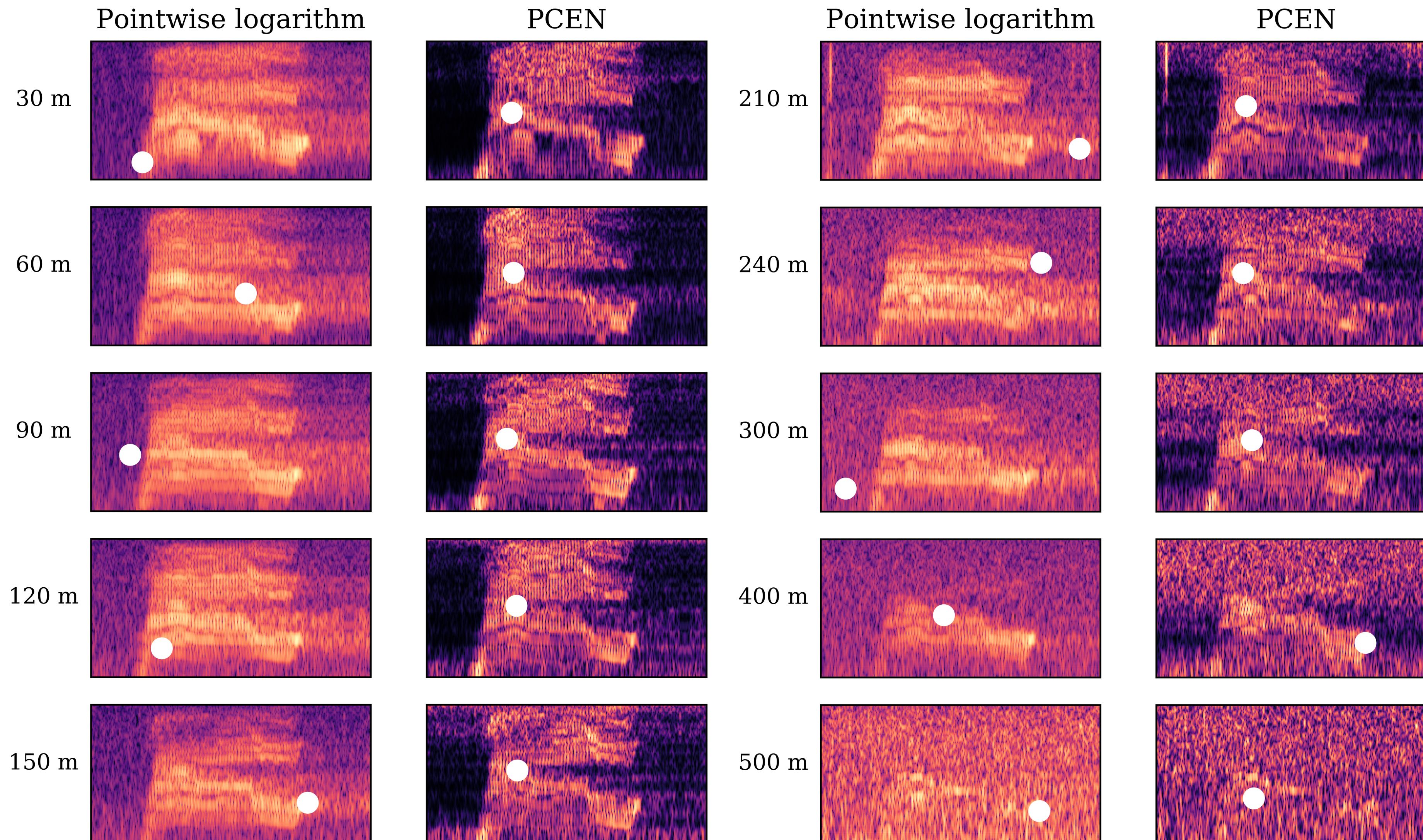
PCEN generalizes (if well-tuned)

[Lostanlen et al. 2019b]



PCEN improves detection range

[Lostanlen et al. 2019c. Audio data from Elly Knight]



End-to-end PCEN learning

[Wang et al. 2017, Zinemanas et al. 2018, Zeghidour et al. 2021, many others]

LEAF (Learnable Audio Frontend) from Google Research.

State-of-the-art performance on many tasks:

speech commands, **birdsong** detection, **musical instruments** ...

Table 2: Impact of the compression layer on the performance (single task, average accuracy in %).

	mel-filterbank			LEAF- filterbank		
	log	PCEN	sPCEN	log	PCEN	sPCEN
Average	73.9 ± 0.8	76.4 ± 0.8	76.0 ± 0.7	74.6 ± 0.7	76.4 ± 0.8	76.9 ± 0.8

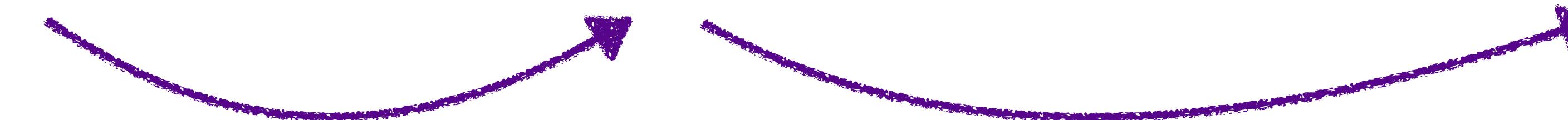


Table from Zeghidour et al. 2021

Conferences > 2019 IEEE Workshop on Applica... 

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

Publisher: IEEE

[Cite This](#)

 PDF

Mark Cartwright ; Jason Cramer ; Justin Salamon ; Juan Pablo Bello [All Authors](#)

Conferences > 2019 IEEE International Smart... 

Activating accessible pedestrian signals by voice using keyword spotting systems

Publisher: IEEE

[Cite This](#)

 PDF

Mirzodaler Muhsinzoda ; Carlos Cruz Corona ; David A. Pelta ; Jose Luis Verdegay [All Authors](#)

[Submitted on 7 Oct 2021]

Improving Bird Classification with Unsupervised Sound Separation

Tom Denton, Scott Wisdom, John R. Hershey

[Submitted on 24 Sep 2021]

Parameterized Channel Normalization for Far-field Deep Speaker Verification

Xuechen Liu, Md Sahidullah, Tomi Kinnunen

Conferences > ICASSP 2021 - 2021 IEEE Inter... 

Sound Event Detection in Urban Audio with Single and Multi-Rate Pcen

Publisher: IEEE

[Cite This](#)

 PDF

Christopher Ick ; Brian McFee [All Authors](#)

Conferences > ICASSP 2020 - 2020 IEEE Inter... 

Audio Feature Extraction for Vehicle Engine Noise Classification

Publisher: IEEE

[Cite This](#)

 PDF

Luca Becker ; Alexandru Nelus ; Johannes Gauer ; Lars Rudolph ; Rainer Martin [All Authors](#)

Conferences > 2019 24th Conference of Open ... 

End-to-end Convolutional Neural Networks for Sound Event Detection in Urban Environments

Publisher: IEEE

[Cite This](#)

 PDF

Pablo Zinemanas ; Pablo Cancela ; Martín Rocamora [All Authors](#)

Conferences > 2019 IEEE Workshop on Applica... [?](#)

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

Publisher: IEEE

[Cite This](#)

[PDF](#)

Mark Car...

Bright ; Jason Cramer ; Justin Salomon ; Juan Pablo Bello [All Authors](#)

conda install -c conda-forge librosa

Activating accessible pedestrian signals by voice using keyword spotting systems

Publisher: IEEE

[Cite This](#)

[PDF](#)

Mirzodaler Muhsinzoda ; Carlos Cruz Corona ; David A. Pelta ; Jose Luis Verdegay [All Authors](#)

[Submitted on 7 Oct 2021]

Improving Bird Classification with Unsupervised Sound Separation

Tom Denton, Scott Wisdom, John R. Hershey

[Submitted on 24 Sep 2021]

Parameterized Channel Normalization for Far-field Deep Speaker Verification

Xuechen Liu, Md Sahidullah, Tomi Kinnunen

Conferences > ICASSP 2021 - 2021 IEEE Inter... [?](#)

Sound Event Detection in Urban Audio with Single and Multi-Rate Pcen

Publisher: IEEE

[Cite This](#)

[PDF](#)

Christopher Ick ; Brian McFee [All Authors](#)

Conferences > ICASSP 2020 - 2020 IEEE Inter... [?](#)

Audio Feature Extraction for Vehicle Engine Noise Classification

Publisher: IEEE

[Cite This](#)

[PDF](#)

Luca Becker ; Alexandru Nelus ; Johannes Gauer ; Lars Rudolph ; Rainer Martin [All Authors](#)

Conferences > 2019 24th Conference of Open ... [?](#)

End-to-end Convolutional Neural Networks for Sound Event Detection in Urban Environments

Publisher: IEEE

[Cite This](#)

[PDF](#)

Pablo Zinemanas ; Pablo Cancela ; Martín Rocamora [All Authors](#)