Microsoft

# Spatial audio feature discovery with convolutional neural networks

Etienne Thuillier, Hannes Gamper, Ivan J. Tashev
Microsoft Research Redmond

# Need

Better understanding of elevation cues
... for efficient & immersive spatial audio rendering

# Problem

Listening test studies
... *do not scale to large sample sizes*
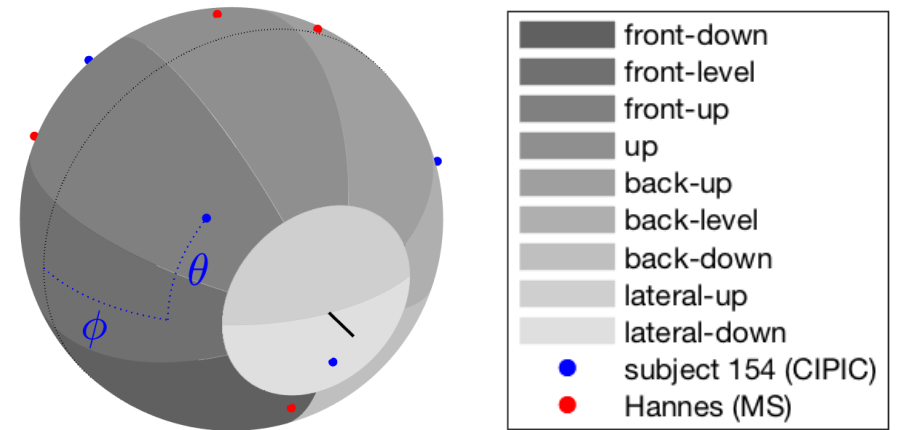... *are exposed to human error*
... *provide highly compressed data*

# Proposed approach (2 steps)
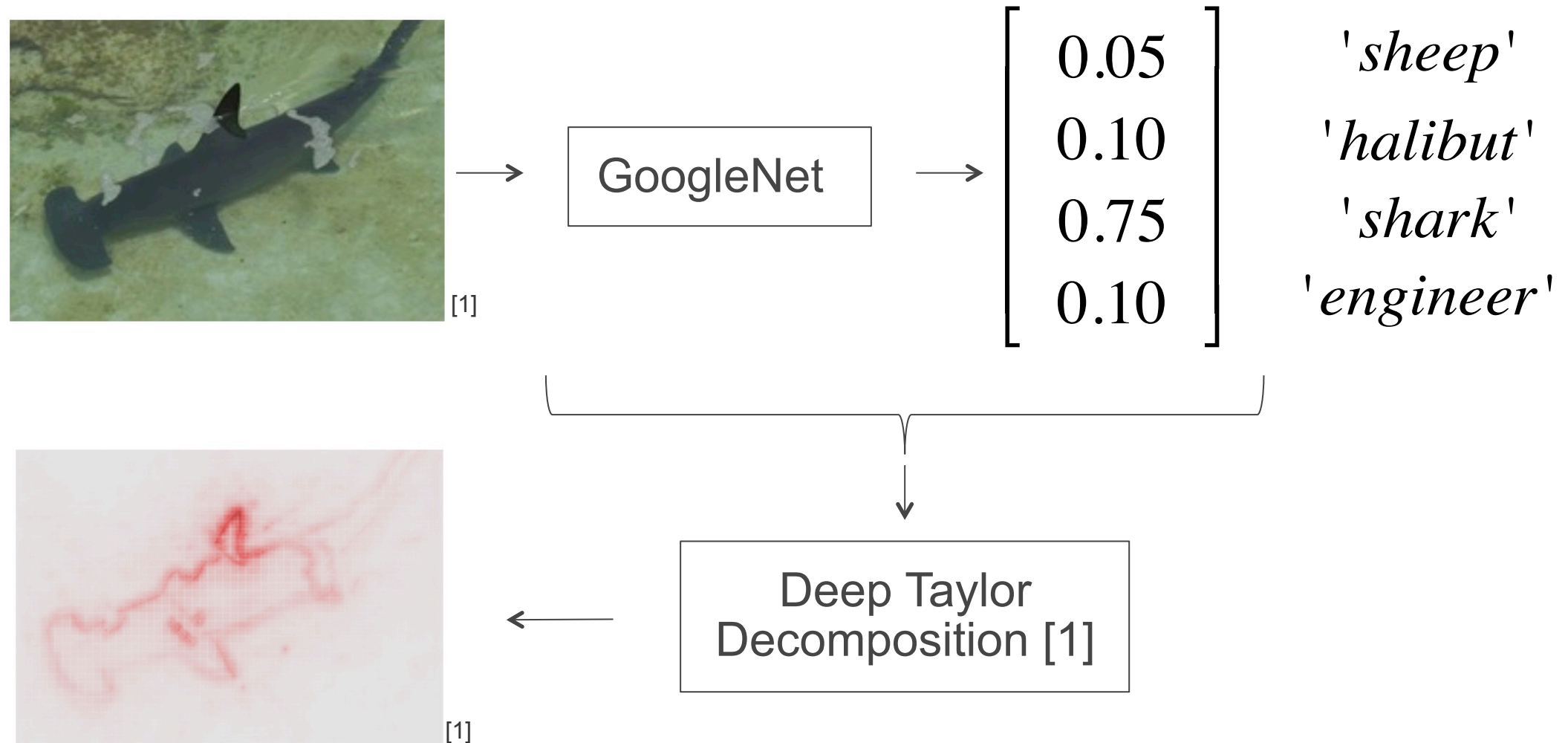
# 1) Train a CNN to classify spatialized sounds

| Dataset | Year | Subjects | Pairs/subj. |
|---|---|---|---|
| Aachen | 2016 | 46 | 2304 |
| Microsoft | 2015 | 252 | 400 |
| RIEC | 2014 | 105 | 865 |
| ARI | 2010 | 135 | 1150 |
| CIPIC | 2001 | 45 | 1250 |

sample HRIR pair

spatialise 50 ms
white noise burst
→



| | |
|---|---|
| | front-down |
| | front-level |
| | front-up |
| | up |
| | back-up |
| | back-level |
| | back-down |
| | lateral-up |
| | lateral-down |
| ● | subject 154 (CIPIC) |
| ● | Hannes (MS) |

~500k HRIR pairs from 583 subjects!

# 2) Apply an explanation technique



$$\begin{bmatrix} 0.05 \\ 0.10 \\ 0.75 \\ 0.10 \end{bmatrix}$$

GoogleNet

*'sheep'*
*'halibut'*
*'shark'*
*'engineer'*

Deep Taylor Decomposition [1]

[1]  G. Montavon et al., "Explaining nonlinear classification decisions with deep  taylor decomposition," *Pattern Recognition*, vol. 65, pp. 211–222, 2017.

# 2) Apply an explanation technique

$$\left[ \ |\mathcal{F}(\boldsymbol{n} * \boldsymbol{h}_\mathrm{i})| \quad |\mathcal{F}(\boldsymbol{n} * \boldsymbol{h}_\mathrm{c})| \ \right] \rightarrow 20\log_{10}(.) \rightarrow$$



ipsi  cont.

kHz

CNN

$$\begin{bmatrix} 0.01 \\ 0.02 \\ 0.75 \\ 0.10 \\ 0.04 \\ 0.03 \\ 0.01 \\ 0.03 \\ 0.01 \end{bmatrix}$$

front-down
front-level
front-up
up
back-up
back-level
back-down
lateral-up
lateral-down

ipsi-lateral        contra-lateral        ipsi  cont.

back-down
back-level
back-up
up
front-up
front-level
front-down

4   8   12        4   8   12
[kHz]              [kHz]

kHz

Deep Taylor
Decomposition

# Models

# Models

| | input | conv. 1 | conv. 2 | conv. 3 | conv. 4 | filters/layer | stride | dense | output layer |
|---|---|---|---|---|---|---|---|---|---|
| **WB** | **> 300 Hz** | $1005 \times 2$ | $\mathbf{25 \times 2}$ | $11 \times 1$ | $11 \times 1$ | $\mathbf{10 \times 1}$ | 4 | 2 | $216 \times 9$ | soft-max |
| **HP** | **> 4 kHz** | $1005 \times 2$ | $\mathbf{25 \times 1}$ | $11 \times 1$ | $11 \times 1$ | $\mathbf{10 \times 2}$ | 4 | 2 | $216 \times 9$ | soft-max |

< 3k trained parameters
> 500 training points parameter
~ 500 subjects in train set

~ 30% classification error (test)

| | CE [%] | RMSE [deg] | MAE [deg] | r |
|---|---|---|---|---|
| random | 91.3 | 74.5 | 59.5 | 0.65 |
| [15] | - | 25.2 | - | 0.85 |
| [27] | - | - | 22.3 | 0.82 |
| [28] | - | - | $\approx 25$ | - |
| **WB** | 45.1 | 43.2 | 16.5 | 0.90 |

1st layer (WB)

S

$\odot$

g( )

1st layer (HP)

S

$\odot$

g( )

$\odot$ element-wise multiply
g( ) sum and rectify
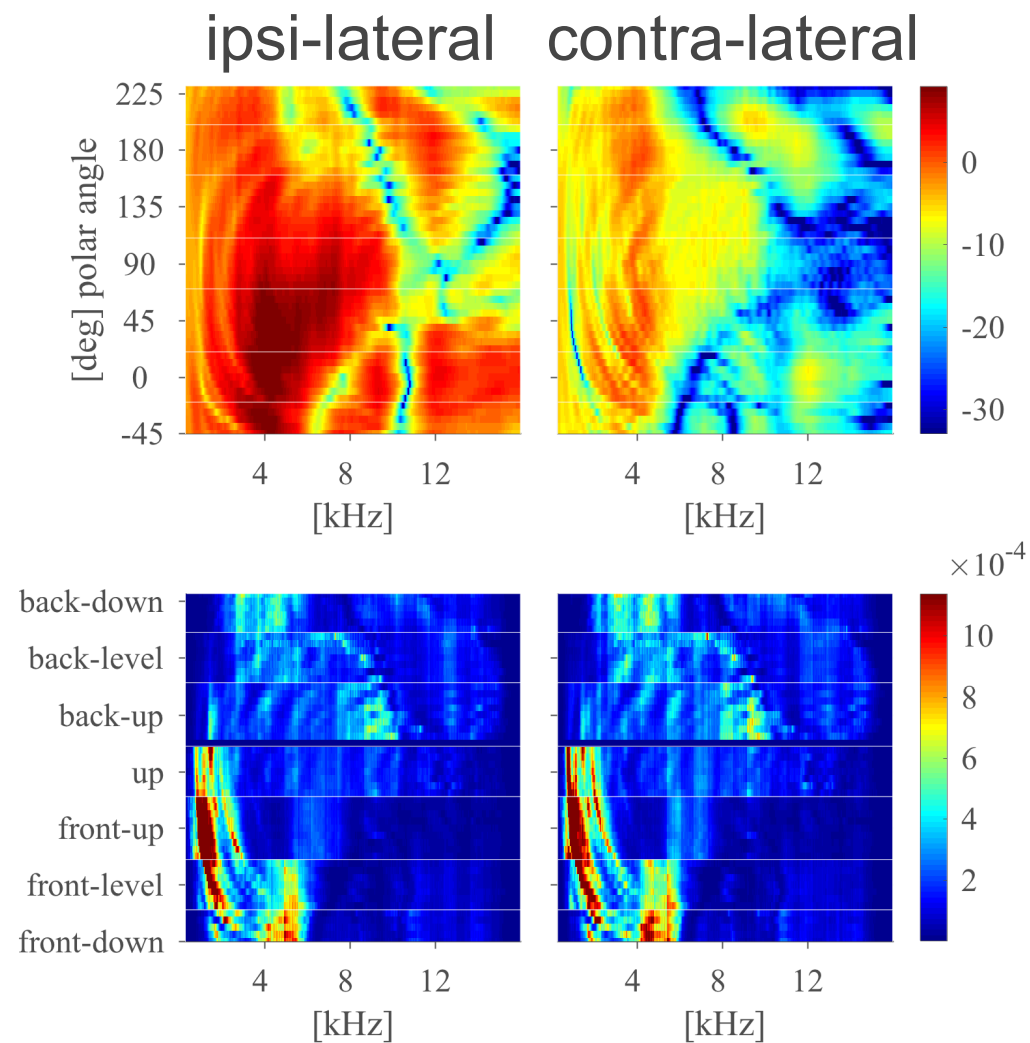
# Results

# HP model

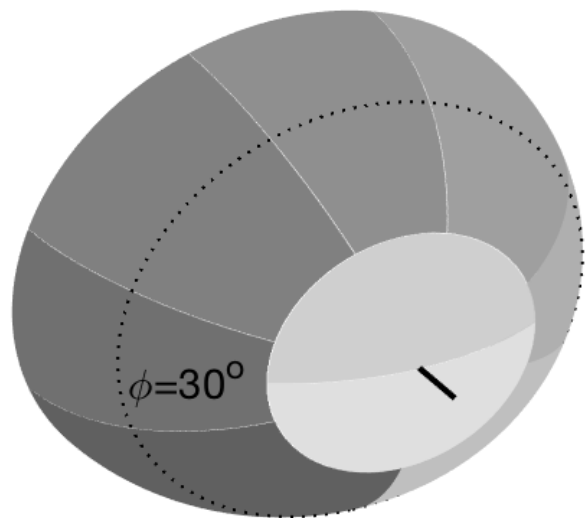Cone of confusion at 30º lat.

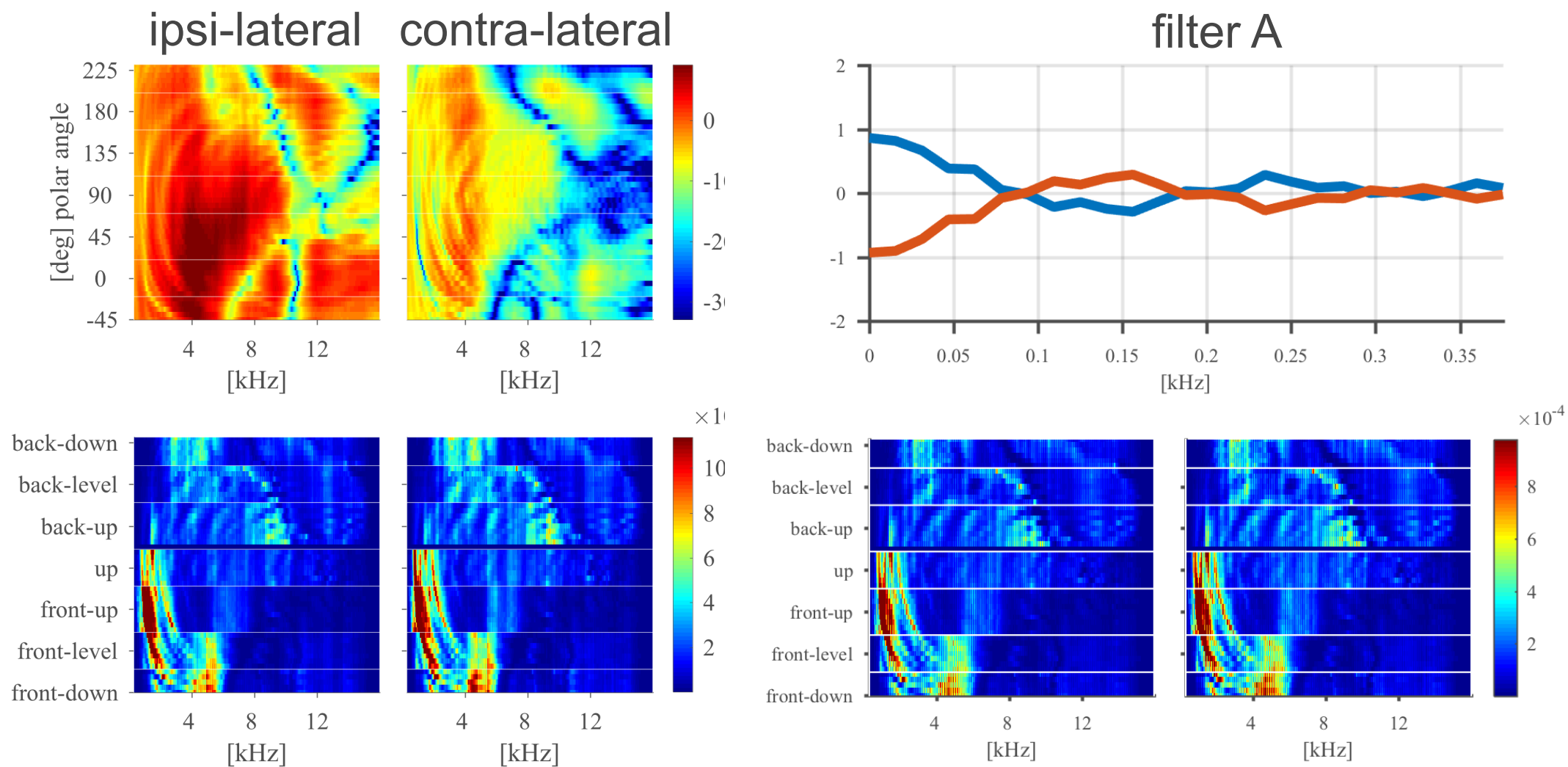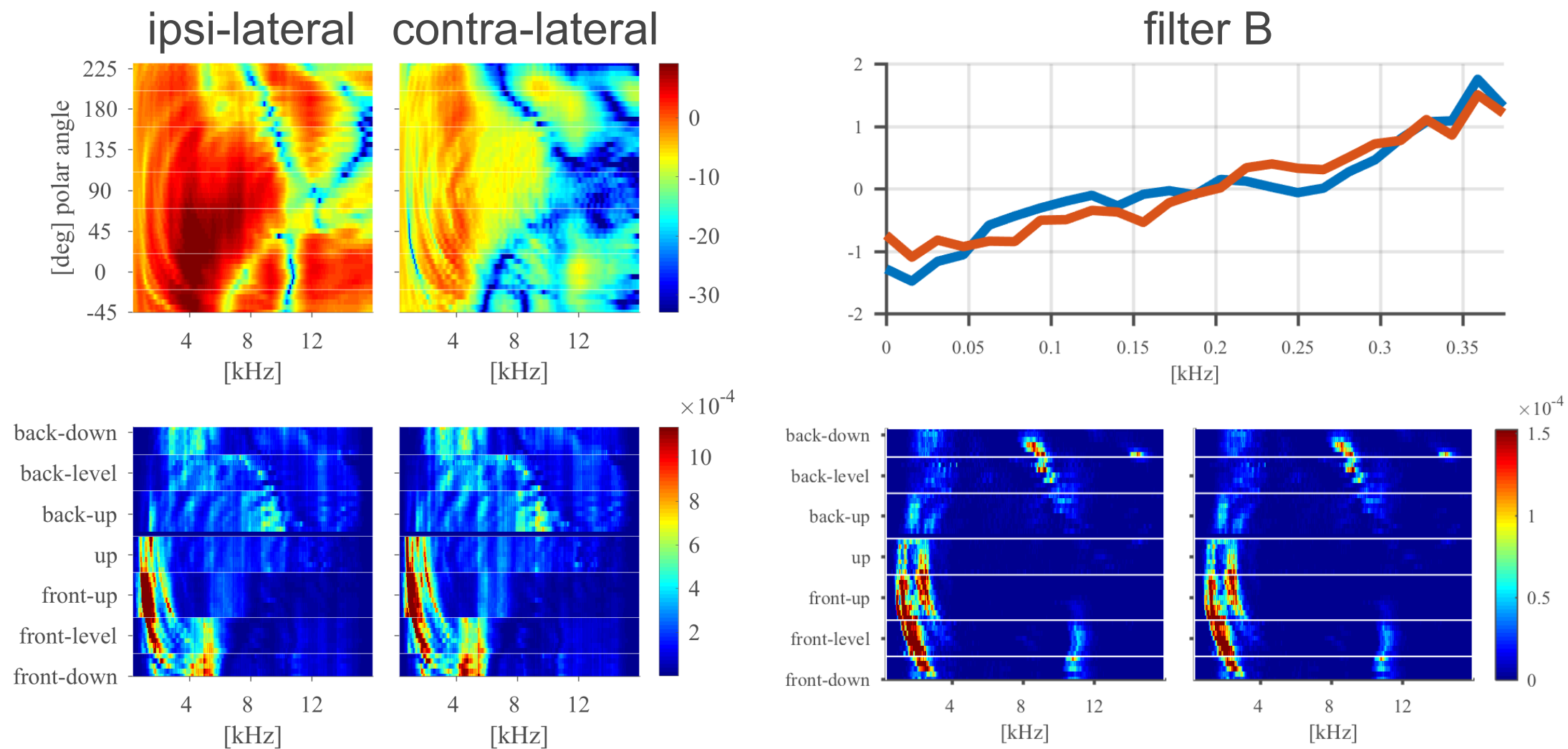CIPIC's subject 154

# WB model

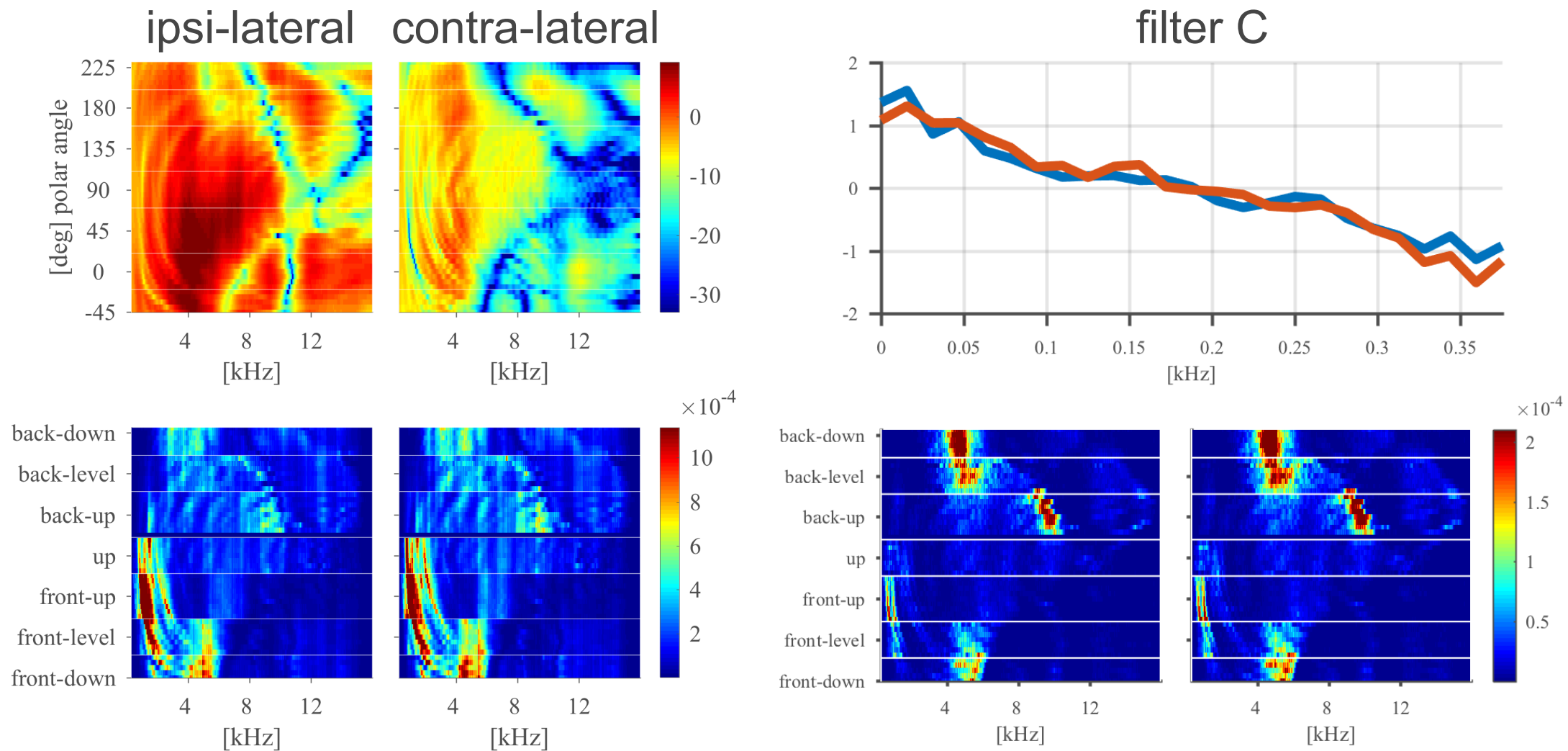Cone of confusion at 30° lat.

CIPIC's subject 154

# WB model

# WB model

# WB model

# Conclusion

# Summary and conclusion

- (Spatial) audio feature discovery using neural network explanation techniques

- Rudimentary models seem to learn cues similar to ones reported in listening test experiments

- Deep Taylor Decomposition seems useful for discovering/visualizing audio features

- Results to be confirmed using a perceptually-grounded front-end

- Lots of future work…!

Thank you!