# **Efficient Audio-based Convolutional Neural Networks via Filter Pruning Arshdeep Singh**<sup>1</sup> & Mark D. Plumbley Here Al for sound



CENTER FOR VISION, SPEECH AND SIGNAL PROCESSING (CVSSP), FACULTY OF ENGINEERING AND PHYSICAL SCIENCES, UNIVERSITY OF SURREY, UK. <sup>1</sup>Research Fellow A, Email: arshdeep.singh@surrey.ac.uk, ORCID: https://orcid.org/0000-0003-3465-0952

### INTRODUCTION Convolutional Neural Network (CNN) $n^{th}$ layer feature maps feature map teature maps Parrot Audio recordin \*\*\*\* Classification Crow layer ₹0 1.5 ↓ 3 Dog Filter A feature map 0 1 0 0 1 (1 imes 3) + (0 imes 4) + (0 imes 2) + (0 imes 3)Input Multipy-accumulate Operations (MACs)

## **OBTAINING AN EFFICIENT CNN BY ELIMINATING FILTERS**

**Hypothesis:** Similar filters produce similar output and hence, mostly contribute to redundancy and can be eliminated.



Selected filter representatives

**Steps:** Unpruned CNN  $\rightarrow$  Eliminate one of the similar filters  $\rightarrow$  Pruned network  $\rightarrow$  Fine-tuning.

### REFERENCES

[1] Q. Kong et al., "PANNs: Large-scale pretrained audio neural networks for audio pattern recognition," IEEE/ACM TALSP, vol. 28, pp. 2880–2894, 2020. [2] Q. Wang et al., "Looking closer at the scene: Multiscale representation learning for remote sensing image scene classification ," IEEE Transactions on NNLS (in press), pp. 1–15, 2020. [3] K. Kahatapitiya and R. Rodrigo, "Exploiting the redundancy in convolutional filters for parameter reduction," proceedings of the IEEE/CVF WACV, pp. 1410–1420, 2021.

[4] Martín-Morató, Irene, et al., "Low- complexity acoustic scene classification for multi-device audio: Analysis of DCASE 2021 challenge systems," DCAEE workshop 85–89, 2021.

- Convolutional neural networks (CNNs) have the capability to learn from examples or experience like humans.
- CNNs have shown state-of-the-art performance in audio classification [1], image scene classification [2] etc.
- However, CNNs are resource hungry due to their large size and heavy computations.
- This makes a bottleneck to deploy CNNs on resourceconstrained devices such as smart phones or IoTs.
- Moreover, CNNs may have redundancy in their parameters or feature maps [3] (See Slide 2).
- Training CNNs for more time generates more CO<sub>2</sub>, e.g. running GTX 1080 Ti hardware for 2 days generates  $CO_2$ = 5.18 kg  $\equiv$  Driving an average car for 20Km.











## PERFORMANCE ANALYSIS

### Experiments are performed on DCASE 2021 Task1A baseline network for audio scene classification (ASC) [4].

- [a] Openly available tool is used to estimate CO<sub>2</sub> emission Link: https://mlco2.github.io/impact/##compute
- [b] Experiments are performed on openly available ASC dataset and CNN [4].
- Proposed code is openly available at Gitlab. Link:https://gitlab.surrey.ac.uk/as0150/passive-pruning

## ACKNOWLEDGEMENT

This work was supported by Engineering and Physical Sciences Research Council (EPSRC) Grant EP/T019751/1 AI for Sound. Thanks to other AI4S team members: Prof. Mark D. Plumbley, Dr. Helen Cooper, Dr Emily Corrigan-Kavanagh and Andres Fernandez for their suggestions and support. AI4S project link: https://ai4s.surrey.ac.uk/index.



### CONCLUSIONS

• Reduced 25% parameters and 27% MACs, with less than 1% drop in accuracy.

• Fine-tuning using few training examples improves the performance of the pruned network significantly.

• Designing better similarity measure and reducing complexity in fine-tuning is a future goal.