

# MONITORING OF DOMESTIC ACTIVITIES BASED ON MULTI-CHANNEL ACOUSTICS A TIME-CHANNEL 2D-CONVOLUTIONAL APPROACH

Technical Report

*Marco Tiraboschi*

Università degli Studi di Milano,  
Department of Computer Science, Laboratorio di Informatica Musicale  
via Celoria 18, Milan, MI 20133, Italy  
marco.tiraboschi@studenti.unimi.it

## ABSTRACT

This approach is meant to be an extension of the DCASE 2018 task 5 baseline system for domestic activity recognition exploiting multi-channel audio: the Convolutional Neural Network model has been slightly restructured for this purpose by using two-dimensional convolutions along the dimensions of time and channel.

## 1. PREPROCESSING

The CNN input layer is the matrix of the short-time log mel-band energies for each one of the audio channels: 500 windows of 40ms with 50% overlap are used for the STFT computation and 40 log mel-band energies are extracted for each time window.

The result is a three-dimensional  $4 \times 500 \times 40$  matrix: 4 audio channels, 500 time windows and 40 log mel-band energies.

## 2. SYSTEM

The network architecture is based on the baseline system convolutional network [1] and extended by operating the convolutions along the two dimensions of time and channel.

The rationale behind this choice is that the network will learn patterns that are equivariant along the dimensions of time and channel, but not along the frequency axis, as similar patterns in different frequency bands don't necessarily belong to similar acoustic events.

- Input data:  $4 \times 500 \times 40$
- Architecture:
  - 2D Convolutional layer (filters: 64, kernel size:  $3 \times 5$ , stride: 1, axes: channel, time) + Batch Normalization + ReLU activation
  - 1D Max Pooling (pool size: 5, stride: 5, axis: time) + Dropout (rate: 20%)
  - 2D Convolutional layer (filters: 64, kernel size:  $2 \times 3$ , stride: 1, axes: channel, time) + Batch Normalization + ReLU activation
  - 1D Global Max Pooling (axis: time) + Dropout (rate: 20%)
  - Dense layer (neurons: 128) + ReLU activation + Dropout (rate: 20%)
  - Softmax output layer (classes: 9)

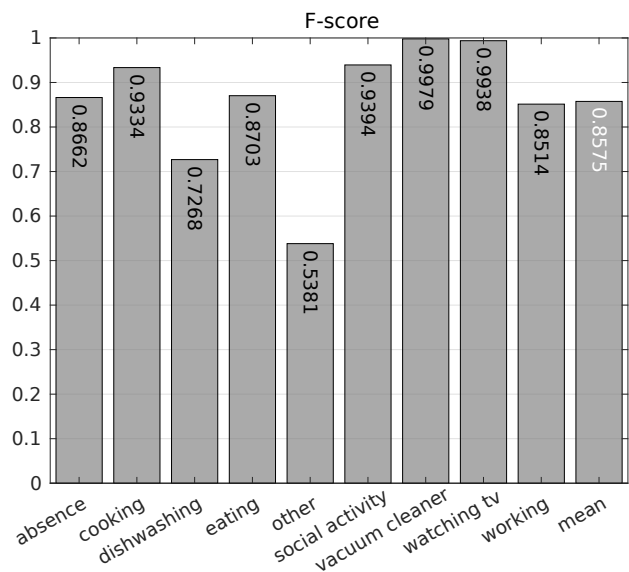


Figure 1: Class-wise and macro-averaged F-scores resulted from cross-validation.

Because of time and computational resources constraints, the model has been trained for 80 epochs only. In the future, it would be interesting to train the model for 420 epochs more to match the baseline training time for a fairer measure of its performance.

- Learning:
  - Optimizer: Adam (learning rate: 0.0001)
  - Epochs: 80
  - Batch size: 512

The proposed system has been implemented in Matlab, using the Neural Network Toolbox [2].

## 3. CROSS-VALIDATION

As output of the training, four convolutional networks have been obtained, each by training the model on a different fold of the cross-validation setup defined by the challenge coordinators. Here are

Output Class	absence	<b>15943</b> 21.8%	<b>1</b> 0.0%	<b>2</b> 0.0%	<b>18</b> 0.0%	<b>344</b> 0.5%	<b>48</b> 0.1%	<b>0</b> 0.0%	<b>4</b> 0.0%	<b>1590</b> 2.2%	<b>88.8%</b> <b>11.2%</b>
	cooking	<b>2</b> 0.0%	<b>4923</b> 6.7%	<b>221</b> 0.3%	<b>19</b> 0.0%	<b>86</b> 0.1%	<b>67</b> 0.1%	<b>3</b> 0.0%	<b>1</b> 0.0%	<b>102</b> 0.1%	<b>90.8%</b> <b>9.2%</b>
	dishwashing	<b>4</b> 0.0%	<b>138</b> 0.2%	<b>1003</b> 1.4%	<b>59</b> 0.1%	<b>59</b> 0.1%	<b>17</b> 0.0%	<b>0</b> 0.0%	<b>0</b> 0.0%	<b>56</b> 0.1%	<b>75.1%</b> <b>24.9%</b>
	eating	<b>9</b> 0.0%	<b>8</b> 0.0%	<b>24</b> 0.0%	<b>1903</b> 2.6%	<b>37</b> 0.1%	<b>31</b> 0.0%	<b>0</b> 0.0%	<b>5</b> 0.0%	<b>48</b> 0.1%	<b>92.2%</b> <b>7.8%</b>
	other	<b>78</b> 0.1%	<b>43</b> 0.1%	<b>94</b> 0.1%	<b>59</b> 0.1%	<b>936</b> 1.3%	<b>40</b> 0.1%	<b>0</b> 0.0%	<b>1</b> 0.0%	<b>168</b> 0.2%	<b>66.0%</b> <b>34.0%</b>
	social activity	<b>12</b> 0.0%	<b>0</b> 0.0%	<b>7</b> 0.0%	<b>5</b> 0.0%	<b>14</b> 0.0%	<b>4479</b> 6.1%	<b>0</b> 0.0%	<b>51</b> 0.1%	<b>24</b> 0.0%	<b>97.5%</b> <b>2.5%</b>
	vacuum cleaner	<b>0</b> 0.0%	<b>1</b> 0.0%	<b>0</b> 0.0%	<b>0</b> 0.0%	<b>0</b> 0.0%	<b>0</b> 0.0%	<b>969</b> 1.3%	<b>0</b> 0.0%	<b>0</b> 0.0%	<b>99.9%</b> <b>0.1%</b>
	watching tv	<b>1</b> 0.0%	<b>0</b> 0.0%	<b>0</b> 0.0%	<b>1</b> 0.0%	<b>0</b> 0.0%	<b>162</b> 0.2%	<b>0</b> 0.0%	<b>18581</b> 25.5%	<b>0</b> 0.0%	<b>99.1%</b> <b>0.9%</b>
	working	<b>2811</b> 3.9%	<b>10</b> 0.0%	<b>73</b> 0.1%	<b>244</b> 0.3%	<b>584</b> 0.8%	<b>100</b> 0.1%	<b>0</b> 0.0%	<b>5</b> 0.0%	<b>16656</b> 22.8%	<b>81.3%</b> <b>18.7%</b>
			<b>84.5%</b> <b>15.5%</b>	<b>96.1%</b> <b>3.9%</b>	<b>70.4%</b> <b>29.6%</b>	<b>82.5%</b> <b>17.5%</b>	<b>45.4%</b> <b>54.6%</b>	<b>90.6%</b> <b>9.4%</b>	<b>99.7%</b> <b>0.3%</b>	<b>99.6%</b> <b>0.4%</b>	<b>89.3%</b> <b>10.7%</b>
	Target Class	absence	cooking	dishwashing	eating	other	social activity	vacuum cleaner	watching tv	working	

Figure 2: Confusion matrix of the cross-validation results.

presented the class-wise and mean F-scores (figure 1) and the confusion matrix (figure 2) resulted from cross-validation.

#### 4. EVALUATION OUTPUT

For each file in the evaluation dataset, the output class is the class that maximizes the sum of the prediction scores from the four convolutional networks.

#### 5. CONCLUSIONS

The results obtained with this approach have shown that exploiting inter-channel acoustic features can significantly improve the performance of a system of acoustic event recognition.

In effect, the proposed system exceeds the baseline by 1.25% macro-averaged F-score (85.75% against 84.50%) in 16% of the training time (80 epochs against 500).

#### 6. REFERENCES

- [1] G. Dekkers, L. Vuegen, T. van Waterschoot, B. Vanrumste, and P. Karsmakers, "DCASE 2018 Challenge - Task 5: Monitoring of domestic activities based on multi-channel acoustics," KU Leuven, Tech. Rep., July 2018
- [2] Matlab Neural Network Toolbox <https://www.mathworks.com/help/nnet>