# SEMI-SUPERVISED NMF-CNN FOR SOUND EVENT DETECTION

## Technical Report

*Teck Kai Chan,[1,2]\*, Cheng Siong Chin[1]*

*Ye Li[2]*

[1]Newcastle University Singapore
Faculty of Science, Agriculture, and Engineering
Singapore 599493
{t.k.chan2, cheng.chin}@newcastle.ac.uk

[2]Xylem Water Solution Singapore Pte Ltd
3A International Business Park
Singapore 609935
ye.li@xyleminc.com

**ABSTRACT**

For the DCASE 2020 Challenge Task 4, this paper proposed a combinative approach using Nonnegative Matrix Factorization (NMF) and Convolutional Neural Network (CNN). The main idea begins with utilizing NMF to approximate strong labels for the weakly labeled data. Subsequently, based on the approximated strongly labeled data, two different CNNs are trained using a semi-supervised framework where one CNN is used for clip-level prediction and the other for frame-level prediction. Using this idea, the best model trained can achieve an event-based F1-score of 45.7% on the validation dataset. Using an ensemble of models, the event-based F1-score can be increased to 48.6%. By comparing with the baseline model, the proposed model outperforms the baseline model by a margin of over 8%.

***Index Terms***— Nonnegative matrix factorization, convolutional neural network, semi-supervised learning, DCASE 2020

## 1. INTRODUCTION

A Sound Event Detection (SED) system can be described as an intelligent system that is capable of not only detecting the types of sound events present in an audio recording but also returning the temporal location of the detected events. Such a system can be useful in several different domains and as compared to a visual detection system, it can be advantageous in several different aspects. Firstly, a SED system is not affected by the degree of illumination. Secondly, occluded objects do not affect detection accuracy. Thirdly, audio recording requires lesser computational resources as compared to an image or video. Finally, some events, such as a car horn, can only be detected by sound [1], [2].

However, for a SED system to achieve maximum performance, there may be a need for a large amount of strongly labeled data where the occurrence of each event with its onset and offset is known with certainty during the model development phase. This can be a limiting factor because such data is usually difficult and time-consuming to collect as it requires repeated listening and adjusting of label time boundaries on a visual interface [3].

As shown in our previous work [4], NMF can be used to approximate strong labels for the weakly labeled data. Thus, as a follow-up work, we proposed to label the weakly labeled data using NMF in a supervised manner. Using the approximated strongly labeled data, we then trained two different CNNs in a semi-supervised framework where one of the models will produce the clip level prediction. In contrast, the other model will produce a frame-level prediction.

Based on such an idea, our best model can achieve an event-based F1-score of 45.7% on the validation dataset. Using an ensemble of models, we can further increase the event-based F1-score to 48.6%. By comparing our models with the baseline model, our models outperformed the baseline model with a margin of over 8%.

The paper is organized as follow, Section 2 provides the information on the proposed methodology, section 3 provides the results and discussion, and finally, the paper ends with a conclusion.

## 2. PROPOSED METHODOLOGY

In this section, information on the proposed methodology will be provided in several different subsections.

### 2.1. Audio Preprocessing and Feature Extraction

As the first step of pre-processing, all audio recordings that were longer or shorter than 10s were first truncated or padded to have an equal length of 10s. The processed recordings were then resampled at 22,050 Hz, and spectrograms were tabulated for each recording using a Fast-Fourier Transform (FFT) window size of 2048 (92 ms)

with a hop length of 345 (15.6 ms). Mel-spectrograms were then tabulated using 64 mel filter banks. Based on such a setting, a tabulated mel spectrogram would have a size of 640 by 64, where 640 represents the number of frames, and 64 represents the number of mel bins. Finally, a logarithm operation was applied to obtain the log mel spectrogram, which will be used as input to the training model.

## 2.2. Approximating Strong Labels Using NMF

NMF is an effective multivariate data decomposition method popularized by Lee and Seung [5]. Given a nonnegative matrix $V$ of size $m \times n$, the objective of NMF is to derive two nonnegative matrices, $W$ of size $m \times r$ and $H$ of size $r \times n$ such that $V$ can be approximated by the linear combination of $W$ and $H$. This can be formally be defined as

$$V \approx WH \tag{1}$$

Where $W$ can be interpreted as the dictionary matrix and $H$ can be interpreted as the activation matrix and $r$ can be interpreted as the number of components.

As shown in our previous work [4], NMF can be used to approximate strong labels for the weakly labeled data. However, the methodology introduced in [4] can induce noise into the training data. With the availability of strongly labeled synthetic data, we proposed to approximate strong labels for weakly labeled data using a supervised approach before the calculation of log mel spectrogram.

The first step is to extract the event template from the mel spectrograms to form a dictionary for different event classes. Since the synthetic sound clip can contain multiple events, temporal masking was applied to the mel spectrogram using the given temporal annotations. Templates of each event class were then retrieved from the masked mel spectrogram using NMF by allowing $r$ to be set as 1. For example, if synthetic clip A has Speech and Cat occurring at frame 1 to 100 and 100 to 110 respectively, all frames from 101 onwards were masked to extract the Speech template followed by masking all frames except frames 100 to 110 to extract the Cat template.

As weakly labeled data possessed the event tags, we then applied the corresponding dictionary on the sound clip to derive $H$. Frames that were activated (above a pre-defined threshold) were assumed to contain the event class. For example, if Clip B contains Speech and Dog, we first applied NMF to decompose Clip B using Speech dictionary and with $r$ set as 1 to derive $H$. Frames that were over a threshold were then assumed to contain only Speech. A similar concept was applied to derive the temporal annotation for Dog by using the Dog dictionary instead of Speech dictionary.

## 2.3. Semi-Supervised Learning

As mentioned in [6], there can be a trade-off in SED performance due to the pooling operation. While the accuracy of clip level detection (also known as audio tagging) can be improved with higher temporal compression (pooling along the time axis), this can result in a degradation of accuracy in frame-level detection.
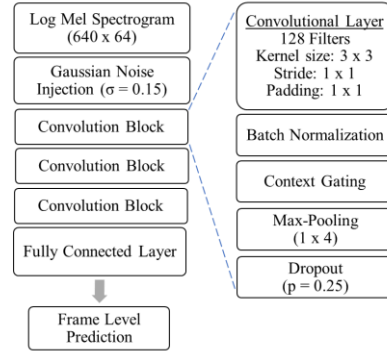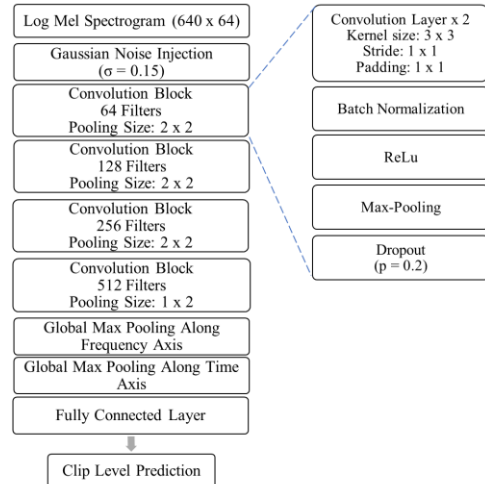


Fig. 1. Model for Frame Level Prediction



Fig. 2. Model for Clip Level Prediction

Therefore, we proposed a Shallow Model (SM) with no temporal compression for frame-level prediction and a Deep Model (DM) with temporal compression for clip level prediction. In addition to the difference in pooling size, SM has fewer convolutional layers, adopted context gating [7] as the activation function as opposed to ReLu and has a slightly higher dropout rate. The details of SM and DM can be found in Fig. 1 and Fig. 2 respectively.

Given that $y_f$ and $y_c$ are the frame level and clip level ground truth of an input respectively. The Binary Cross Entropy (BCE) loss between the frame-level prediction of SM, $SM_f$, and $y_f$ can be given as

$$l_f = BCE(SM_f, y_f) \tag{2}$$

While the BCE loss between the clip level prediction of DM, $DM_c$, and $y_c$ can be given as

$$l_c = BCE(DM_c, y_c) \qquad (3)$$

As mentioned earlier, the accuracy of clip-level detection is better for models with higher temporal compression. We hypothesized that by enforcing the prediction of SM to be consistent with DM, it could produce a better frame-level prediction. As the prediction output of SM is in frame level, we applied a global max pooling on the time axis of $SM_f$ to obtain the clip level prediction, $SM_c$. Instead of using BCE as the function, we proposed the use of Mean Square Error (MSE) as the consistency loss function. This was found to be a better consistency loss function as compared to using BCE [8].

$$l_{con} = MSE(SM_c, DM_c) \qquad (4)$$

In addition, we also enforce the consistency of prediction on the unlabeled data. This is also known as semi-supervised learning, which was found to improve the performance and generalization of the model [9], [10]. Given that the frame level and clip level prediction from SM and DM on the unlabeled data are $SM_{uf}$ and $DM_{uc}$ respectively. The clip level prediction, $SM_{uc}$, from SM can be obtained by applying a global max pooling on $SM_{uf}$. Thus the consistency cost on the unlabeled data can be given as

$$l_{unlabel} = MSE(SM_{uc}, DM_{uc}) \qquad (5)$$

However, if $l_{unlabel}$ was given too much weightage during the early stage of training, it may result in a degenerate solution where no meaningful classification of the data can be obtained [8]. Thus, a weighting parameter, $w$, is required to regularize the contribution of $l_{unlabel}$ throughout the entire training process. Following [8], $w$ was proposed to ramp up from 0 along a Gaussian curve and can be defined as

$$w = \exp\left(-5(1-T)^2\right) \qquad (6)$$

Where $T$ is a positive value which represents the training progression. As an additional measure to prevent obtaining a degenerate solution, we proposed allowing $l_{unlabel}$ to be calculated if DM is confident with its prediction. Thus $l_{unlabel}$ is defined as

$$l_{unlabel} = \begin{cases} w \times MSE(SM_{uc}, DM_{uc}), \max(DM_{uc}) > \lambda \\ 0, otherwise \end{cases} \qquad (7)$$

Where $\lambda$ represent the level of confidence. Thus, the total combined loss, $l_{total}$ can then be defined as

$$l_{total} = l_f + l_c + l_{con} + l_{unlabel} \qquad (8)$$

Based on the calculated $l_{total}$, model parameters of both models will then be updated using Adam with its default parameters [11]. As it was found that the performance of deep NN may benefit from resetting the Learn-ing Rate (LR) during the training process [12], we proposed to anneal the LR according to a cosine function and reset it to original LR after a certain number of epochs. The LR at each iteration is defined as [12]

$$LR_{curr} = LR_{max} + \frac{1}{2}\left(LR_{max} - LR_{min}\right)\left(1 + \cos\left(\frac{T_{curr}}{T_i}\pi\right)\right) \qquad (9)$$

Where $LR_{max}$ represents the maximum LR and was set as 0.0012. $LR_{min}$ represents the minimum LR which was set as 1e-6. $T_{curr}$ represents the current training iteration and $T_i$ represent the maximum training iterations before a LR reset. If $T_{curr}$ is equal to 0, the LR will be at its maximum value. If $T_{curr}$ is equal to $T_i$, the LR rate will be at its minimum and at the next iteration, $T_{curr}$ will then be reset to 0 while $T_i$ is multiplied with an integer, $T_{mult}$ which can delay the next restart if $T_{mult}$ is larger than 1.

As $T$ represents the training progression, which directly affects the calculation of $l_{unlabel}$. We proposed to define $T$ as

$$T = \frac{T_{curr}}{T_i} \qquad (10)$$

Thus the contribution $w$ will be reset to 0 whenever the LR is reset.

Finally, we also adopted the concept of transfer learning in our system, where the models will first be trained using synthetic data for 5 epochs without the inclusion of $l_{unlabel}$. Only from the 6th epoch onwards, the parameters will be updated using real data and with the inclusion of $l_{unlabel}$.

## 2.4. Post Processing

A clip is considered to contain a specific event if the predicted probability from DM is larger than 0.5. Using the identified audio tag, temporal location can be found by locating the activated frames based on the predicted outputs from SM.

Before locating the activated frames, we smoothed the outputs from SM using iterative median filter [13] with an event-specific window size. Frames were then considered to be activated if they exceeded an event-specific frame threshold.

Following the implementation in [14], neighboring frames were also considered to be activated if they exceeded a lower bound threshold of 0.08. In addition, detected events with a duration of shorter than 0.1s were removed as they were considered as noise. Finally, we concatenated two similar events together if the difference between the first event offset and the second event onset is shorter than 0.2s.

| | Event-Based F1-Score (%) | PSDS F-Score (%) | PSDS | PSDS Cross Trigger | PSDS Macro |
|---|---|---|---|---|---|
| PS 1 ($T_i$ = 1 epochs, $T_{mult}$ = 2, $LR_{min}$ = 1e-6) * | 45.2 | 63.6 | 0.630 | 0.548 | 0.408 |
| PS 2 ($T_i$ = 1 epochs, $T_{mult}$ = 2, $LR_{min}$ = 1e-6, trained with Synthetic Data) * | 45.7 | 65.6 | 0.635 | 0.546 | 0.409 |
| Ensemble System (PS 1 + PS 2) * | 48.0 | 66.6 | 0.652 | 0.577 | 0.430 |
| Ensemble System (PS 1 + PS 2) with Tuned MF Window * | 48.6 | 66.5 | 0.649 | 0.573 | 0.425 |
| Baseline without Source Separation | 34.8 | 60.0 | 0.61 | 0.524 | 0.433 |
| Baseline with Source Separation | 35.6 | 60.5 | 0.626 | 0.546 | 0.449 |

Table 2. Results of Proposed Methodology Against Baseline Systems (PS refers to Proposed System and system with * were the submitted system to the DCASE Challenge Task 4)

## 3. RESULTS AND DISCUSSION

In our experiment, several factors can influence the training process and eventually affect the detection accuracy. Firstly, we found that if $\lambda$ was set as a low value (i.e. 0.5), this will produce a higher $l_{unlabel}$, which resulted in a suboptimal solution. Thus, $\lambda$ was set as 0.9 to ensure that $l_{unlabel}$ will only be calculated based on highly confident prediction.

| Event | Frame Threshold | First MF Window Size | Second MF Window Size |
|---|---|---|---|
| Speech | 0.3 | 7 | 15 |
| Dog | 0.3 | 7 | 15 |
| Cat | 0.3 | 3 | 6 |
| Alarm/Bell Ringing | 0.4 | 8 | 21 |
| Dishes | 0.2 | 3 | 5 |
| Frying | 0.6 | 24 | 48 |
| Blender | 0.6 | 3 | 6 |
| Running Water | 0.6 | 4 | 13 |
| Vacuum Cleaner | 0.4 | 24 | 48 |
| Electric Shaver/ Toothbrush | 0.4 | 48 | 96 |

Table. 1. Optimal Global Event Specific Threshold and MF Window Size (in Frames)

Secondly, by using an event-specific frame threshold, detection accuracy can be raised. However, optimal values differ across systems. Likewise, the median filter window is also dependent on the system trained. In our experiments, we found that using a smaller window in the first round of filtering and larger window size in the second round of filtering usually produces higher detection accuracy. Through multiple experiments, we derived the optimal global value for the event-specific frame threshold and filter window, and these are shown in Table 1.

In [14], the post-processing method was to join similar events together before the removal of noise. However, we realized that the accuracy could be higher if the noise were removed before the concatenation of similar events.

$T_i$ is a hyperparameter that controls how fast LR will reduce from $LR_{max}$ to $LR_{min}$. In our experiment, if $T_i$ was a small value (i.e. smaller than 5 epochs), $T_{mult}$ must be at least 2 to prevent the large fluctuation of LR throughout the training process. Whereas if $T_i$ was a large value (more than 5 epochs), $T_{mult}$ can be set as 1 as the transition of LR from $LR_{max}$ to $LR_{min}$ can be considered slow and steady. Subsequently, we found that it is not a guarantee that a better solution can be found following a LR reset, and there is a possibility that a worse solution is found.

For transfer learning, we tested a different number of epochs for the transition of training data from synthetic data to real data. Based on our results, changing the data type after 5 epochs appeared to yield the best results.

Based on our methodology described in Section 2 and the aforementioned findings, our system can achieved an event-based F1-score of 45.2% by allowing $T_i$ to be set as 1 epoch and $T_{mult}$ as 2.

As mentioned earlier, models were trained using only synthetic data for the first 5 epochs and only from the 6th epoch onwards, model were trained using real data. In our experiment, we also tested a different form of transition where model were trained using synthetic data for the first 5 epoch whereas from 6th epoch onwards, models were trained using both real and synthetic data. Based on such setting, our system can achieved a slightly higher event based F1-score of 45.7% by allowing $T_i$ to be set as 1 epoch and $T_{mult}$ as 2.

Using an ensemble of the two models, we can further increase the event-based F1-score to 48.0%. As mentioned earlier, each system can have a different optimal MF window, we then tuned the MF window of the ensemble model, and this could further increase the event-based F1-score to 48.6%. However, we suspect this can have an adverse effect on the system, which may overfit the system to the validation dataset.

Based on the results shown in Table 2, all of the models trained using our proposed methodology were able to win the baseline system by a margin of at least 8%. We also computed the Polyphonic Sound Detection Score

(PSDS) [16] as a secondary measure. While we have a higher event-based F1-score, our system has a lower PSDS Macro score. It could be due to the difficulty of detecting Dishes and had a much lower detection accuracy (<20% across the systems) as compared to the other event class.

## 4. CONCLUSION

In this paper, a combinative approach using NMF and CNN was proposed for DCASE Challenge 2020 task 4. The proposed system could achieve an event-based F1-score of 45.7%, and with the use of the ensemble method, event-based F1-score raised to 48.6%. Based on such results, our system could outperform the baseline by a margin of over 8%. For our future work, we will investigate the cause of low detection accuracy for Dishes and improve our system in this aspect.

## 5. ACKNOWLEDGMENT

We would like to thank Nicolas Turpault for his technical support and also Kong Qiuqiang and Lin Liwei for their prompt replies in regards to our questions on their system. Finally, we would also like to express our gratitude to Giacomo Ferroni for his explanation on the PSDS metric.

## 6. REFERENCES

[1] Q. Kong, Y. Xu, I. Sobieraj, W. Wang, and M. D. Plumbley, "Sound Event Detection and Time–Frequency Segmentation from Weakly Labelled Data," *IEEE/ACM Trans on Audio, Speech, and Language Process.*, vol. 27, no. 4, pp. 777-787, Apr. 2019.

[2] Q. Zhou, Z. Feng and E. Benetos "Adaptive Noise Reduction for Sound Event Detection Using Subband-Weighted NMF," *Sensors*, vol. 19, no. 3206, pp. 1-19, Jul. 2019.

[3] B. Kim and B. Pardo, "Sound Event Detection Using Point-Labeled Data," *2019 IEEE Workshop on Appl. of Signal Process. to Audio and Acoustics*, New Paltz, New York, USA, Oct., 2019, pp. 1-5.

[4] T. K. Chan, C. S. Chin and Y. Li, "Nonnegative Matrix Factorization-Convolutional Neural Network (NMF-CNN) For Sound Event Detection," in *Proc. Detection and Classification of Acoustic Scenes and Events 2019 Workshop (DCASE2019)*, New York, USA, Oct. 2019, pp. 40-44.

[5] D. D. Lee, and H. S. Seung, "Learning the parts of objects by nonnegative matrix factorization," *Nature*, vol. 401, no. 6755, pp. 788-791, Oct. 1999.

[6] L. Lin, X. Wang, H. Liu and Y. Qian, "Guided Learning Convolution System For DCASE 2019 Task 4," in *Proc. Detection and Classification of Acoustic Scenes and Events 2019 Workshop*, New York, USA, Oct. 2019, pp. 134-138.

[7] A. Miech, I. Laptev and J. Sivic, "Learnable pooling with Context Gating for video classification," *arXiv preprint arXiv: 1706.06905*, pp. 1-8, Mar. 2018.

[8] S. Laine and T. Aila, "Temporal Ensembling For Semi-Supervised Learning," in *Proc. 5th Int. Conf. Learning Representations (ICLR 2017)*, Toulon, France, Apr. 2017, pp. 1-13.

[9] A. Oliver, A. Odena, C. Raffel, E. D. Cubuk and I. J. Goodfellow, "Realistic Evaluation of Deep Semi-Supervised Learning Algorithms," in *Proc. 32nd Conf. on Neural Information Process.* Syst. (NeurIPS 2018), Montreal, Canada, Dec. 2018, pp. 1-12.

[10] D. H. Lee, "Pseudo-Label : The Simple and Efficient Semi-Supervised Learning Method for Deep Neural Networks," in *Int. Conf. Mach. Learning 2013 Workshop : Challenges in Representation Learning (WREPL)*, Georgia, USA, Jun. 2013.

[11] D. P. Kingma and J. L. Ba, "ADAM: A Method For Stochastic Optimization," in *Proc. 3rd Int. Conf. Learning Representations (ICLR 2015)*, San Diego, USA, May 2015, pp. 1-15.

[12] I. Loshchilov and F. Hutter, "SGDR: Stochastic Gradient Descent With Warm Restarts," in *Proc. 5th Int. Conf. Learning Representations (ICLR 2017)*, Toulon, France, Apr. 2017, pp. 1-16.

[13] E. A. Castro and D. L. Donoho, "Does Median Filtering Truly Preserve Edges Better Than Linear Filtering?," *The Annals of Statistics*, vol. 37, no. 3, pp. 1172–1206, Apr. 2009.

[14] Q. Kong, Y. Cao, T. Iqbal, Yong Xu, W. Wang, and M. D. Plumbley, "Cross-task learning for audio-tagging, sound event detection spatial localization: DCASE 2019 baseline systems," *arXiv preprint arXiv: 1904.03476*, pp. 1-5.

[15] A. Mesaros, T. Heittola and T. Virtanen, "Metrics for Polyphonic Sound Event Detection," *Applied Sciences*, vol. 6, no. 162, pp. 1-17, May 2016.

[16] C. Bilen, G. Ferroni, F. Tuveri, J. Azcarreta and S. Krstulovic, "A Framework For The Robust Evaluation of Sound Event Detection," in *2020 IEEE Int. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, Barcelona, Spain, May 2020, pp. 61-65.