

# GUIDED LEARNING CONVOLUTION SYSTEM FOR DCASE 2019 TASK 4

Liwei Lin<sup>1,2</sup>, Xiangdong Wang<sup>1</sup>, Hong Liu<sup>1</sup>, YueLiang Qian<sup>1</sup>,

<sup>1</sup>Beijing Key Laboratory of Mobile Computing and Pervasive Device,  
Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China

<sup>2</sup>University of Chinese Academy of Sciences, Beijing, China  
{linliwei17g, xdwang, hliu, ylqian}@ict.ac.cn

## ABSTRACT

In this paper, we describe in detail the system we submitted to DCASE2019 task 4: sound event detection (SED) in domestic environments. We employ a convolutional neural network (CNN) with an embedding-level attention pooling module to solve it. By considering the interference caused by the co-occurrence of multiple events in the unbalanced dataset, we utilize the disentangled feature to raise the performance of the model. To take advantage of the unlabeled data, we adopt Guided Learning for semi-supervised learning. A group of median filters with adaptive window sizes is utilized in the post-processing of output probabilities of the model. We also analyze the effect of the synthetic data on the performance of the model and finally achieve an event-based F-measure of 45.43% on the validation set and an event-based F-measure of 42.7% on the test set. The system we submitted to the challenge achieves the best performance compared to those of other participants.

**Index Terms**— Sound event detection, weakly supervised learning, semi-supervised learning, attention, Guided Learning, Disentangled Feature

## 1. INTRODUCTION

DCASE2019 task 4 is the follow-up to DCASE2018 task 4 [1], which aims at exploring the possibility of the large-scale sound event detection using weakly labeled data (without timestamps) and unlabeled data. Different from DCASE2018 task 4, DCASE2019 task 4 introduces an additional strongly annotated synthetic training set.

Sound event detection (SED) consists in recognizing the presence of sound events in the segment of audio and detecting their onset as well as offset. Due to the high cost of manually labeling data, it is essential to efficiently utilize weakly-labeled data and unlabeled data. Simultaneously, the different physical characteristics of events (such as different duration) and the unbalance of the available training set also increase the difficulty of the multi-class SED in domestic environments. For DCASE2019 task4, there are 5 issues to be resolved:

- 1) How to learn efficiently with weakly-labeled data?
- 2) How to learn efficiently with unbalanced training set?
- 3) How to combine weakly-supervised learning with semi-supervised learning efficiently using weakly-labeled data and unlabeled data?
- 4) How to design a better post-processing method on the output probabilities of the model to detect more accurate boundaries according to the characteristics of each event category?

- 5) Does the strongly annotated synthetic training set help?

In this paper, we present a system to solve all these five issues. For issue 1 and 2, we utilize convolutional neural network (CNN) with the embedding-level attention pooling module and disentangled feature [2] to solve them. For issue 3, we adopt a semi-supervised learning method named Guided Learning [3]. For issue 4, according to varied duration of different event categories, we employ a group of median filters with adaptive window sizes in the post-processing of output probabilities of the model. For issue 5, we simply regard the strongly annotated synthetic training set as a weakly annotated training set and conduct a series of ablation experiments to explore its effects on weakly-supervised learning and unsupervised learning separately.

In the rest of this paper, we introduce our methods in Section 2, describe in detail our experiments in Section 3 and draw conclusions in Section 4.

## 2. METHODS

In this section, we discuss the solution for issue 1 in Section 2.1, the solution for issue 2 in Section 2.2, the solution for issue 3 in Section 2.3 and the solution for issue 4 in Section 2.4.

### 2.1. A CNN model with the embedding-level attention pooling module

In this section, we describe in detail the model we employ. As shown in Figure 1a, the model comprises 3 parts: a feature encoder, an embedding-level attention pooling module and a classifier. The feature encoder encodes the input feature of an audio clip into high-level feature representations. Assuming that there are  $C$  event categories to detect, then the embedding-level attention pooling module integrates these high-level feature representations into  $C$  contextual representations. Eventually, the clip-level probabilities can be obtained by passing this  $C$  contextual representations through the classifier.

As shown in Figure 1b, the feature encoder we employs is composed of a Batch normalization layer [4], 3 Max pooling layers and 3 CNN blocks, each of which consists of a CNN layer, a Batch normalization layer and a ReLU activation layer as shown in Figure 1c. And the classifier for each contextual representation is a fully-connected layer with a Sigmoid activation layer.

The ability of this model to carry out weakly-supervised learning attributes to its embedding-level attention pooling module. Let  $\mathbf{x} = \{x_1, \dots, x_T\}$  be the high-level feature representations generated by the feature encoder and  $\mathbf{y} = \{y_1, \dots, y_C\}$  ( $y_c \in \{0, 1\}$ ) be the groundtruth, where  $T$  denotes the number of total frames of the high-level feature representations.

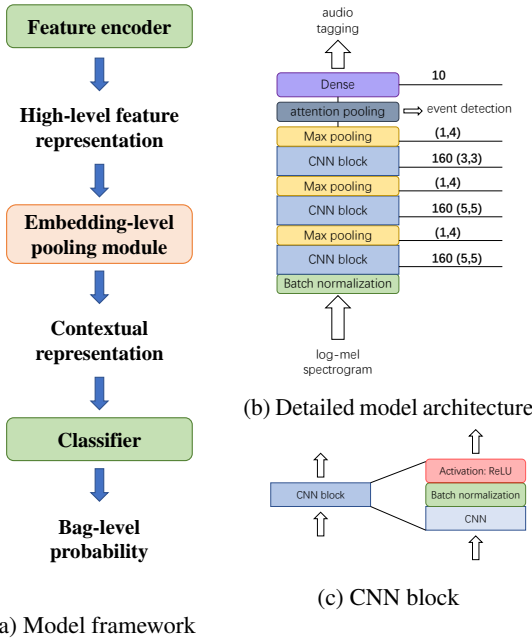


Figure 1: CNN model with the embedding-level attention pooling module.

Then for each category  $c$ , the embedding-level attention pooling gives different weights  $\mathbf{a}_c = \{a_{c1}, \dots, a_{cT}\}$  to the corresponding  $x_t$  in  $\mathbf{x}$ . Then the contextual representation  $\mathbf{h} = \{h_1, h_2, \dots, h_C\}$  can be obtained by the following way:

$$h_c = \sum_t a_{ct} \cdot x_t \quad (1)$$

Such an  $\mathbf{a}_c$  enables the model to treat each frame differently. Important frame  $x_t$  in  $\mathbf{x}$  with larger  $a_{ct}$  contributes more to  $h_c$ . The embedding-level attention pooling module generates  $\mathbf{a}_c$  by the following way:

$$a_{ct} = \frac{\exp((w_c^T x_t + b_c)/d)}{\sum_k \exp((w_c^T x_k + b_c)/d)} \quad (2)$$

where  $d$  is equal with the dimension of  $\mathbf{x}$ ,  $w_c^T$  is a trainable vector, and  $b_c$  is the trainable bias.

More importantly,  $a_{ct}$  possess the ability to indicate key frames of an audio and is able to generate frame-level probabilities as explained in [2]:

$$\hat{p}(y_c | x_t) = \sigma(w_c^T x_t + b_c) \quad (3)$$

where  $\sigma$  is Sigmoid function.

Assuming that  $\hat{\mathbf{P}}(y_c | \mathbf{x})$  is the clip-level probabilities for event category  $c$ , then the clip-level prediction is:

$$\phi_c(\mathbf{x}) = \begin{cases} 1, & \hat{\mathbf{P}}(1 | \mathbf{x}) \geq \alpha \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

The frame-level prediction is:

$$\varphi_c(\mathbf{x}, t) = \begin{cases} 1, & \hat{p}(1 | x_t) \cdot \phi_c(\mathbf{x}) \geq \alpha \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

Without loss of generality, we set  $\alpha = 0.5$ .

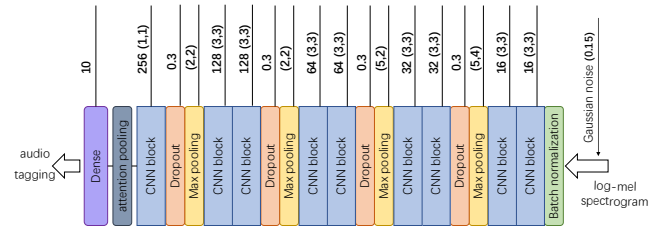


Figure 2: The model architecture of the PT-model.

## 2.2. Disentangled feature

We take disentangled feature (DF) [2], which re-models the high-level feature subspace of each event category according to the prior information without pre-training, to mitigate the effect of the interference caused by the co-occurrence multiple events.

Assuming that  $\chi^d (\mathbf{x} \in \chi^d)$  is a  $d$ -dimensional space generated by the feature encoder and  $\beta = \{e_1, e_2, \dots, e_d\}$  is an orthogonal basis of  $\chi^d$  where the element of  $e_i$  in  $i^{\text{th}}$  dimensional is 1. DF selects specific bases of  $\chi^d$  to construct a specific subspace for each category and the basis of the re-modeled feature space  $\chi_c$  of category  $c$  is

$$\beta'_c = \{e_1, e_2, \dots, e_{k_c}\} \quad (6)$$

$$k_c = \lceil ((1 - m) \cdot f_c + m) \cdot d \rceil \quad (7)$$

$$f_c = \sum_i^C \frac{r_i \cdot N_{ci}}{R} \quad (8)$$

$$R = \max_c \sum_{i=1}^C r_i \cdot N_{ci} \quad (9)$$

where  $m$  is a constant to avoid too-small  $k_c$  and  $N_{ci}$  is the number of clips containing  $i$  categories including category  $c$  in the training set. The constant coefficient  $r_i$  denotes the importance these clips:

$$r_i = \begin{cases} 1, & i = 1 \\ 0, & \text{otherwise} \end{cases} \quad (10)$$

## 2.3. Guided Learning

To combine weakly-supervised learning with semi-supervised learning, we utilize Guide Learning (GL) proposed in [3] with a more professional teacher model (PT-model) to guide a more promising student model (PS-model).

The architecture of the PS-model is consistent with the model described in Section 2.1 and we show that of the PT-model in Figure 2. The CNN feature encoder of the PT-model is considered to be better designed than the PS-model on the audio tagging performance with larger sequential sampling size and less trainable parameters. This is because that the larger sequential sampling size allows the CNN feature encoder of the PT-model to have a larger receptive field followed by better exploitation of contextual information.

However, the larger sequential sampling size also disables the PT-model to see finer information due to the compress of sequential information. Therefore, the PS-model is designed with smaller sequential sampling size to see finer information and then achieves better frame-level prediction.

This gap between their ability makes it possible to optimize the PS-model with the guide of the PT-model using unlabeled data. As

---

**Algorithm 1** Guided learning pseudocode.
 

---

**Require:**  $x_k$  = training input with index  $k$   
**Require:**  $L$  = set of weakly-labeled training input  
**Require:**  $U$  = set of unlabeled training input  
**Require:**  $y_k$  = label of weakly-labeled input  $x_k \in L$   
**Require:**  $S_\theta(x)$  = neural network of the PS-model with trainable parameters  $\theta$   
**Require:**  $T_{\theta'}(x)$  = neural network of the PT-model model with trainable parameters  $\theta'$   
**Require:**  $g(x)$  = stochastic input augmentation function  
**Require:**  $J(t, z)$  = loss function  
**Require:**  $\phi(z)$  = prediction generation function  
**Ensure:**  $\theta, \theta'$

```

for  $i = 1 \rightarrow num\_epochs$  do
  if  $i > start\_epoch$  then
     $a \leftarrow 1 - \gamma^{i-start\_epoch}$   $\triangleright$  calculate the weight of
    unsupervised loss of the PT-model
  else
     $a \leftarrow 0$ 
  end if
  for each minibatch  $\beta$  do
     $s_k \leftarrow S_\theta(x_k \in \beta)$   $\triangleright$  the coarse-level predicted probability
    of the PS-model
     $t_k \leftarrow T_{\theta'}(g(x_k) \in \beta)$   $\triangleright$  the coarse-level predicted
    probability of the PT-model
     $\tilde{s}_k \leftarrow \phi(s_k)$   $\triangleright$  convert the predicted probability into 0-1
    prediction
     $\tilde{t}_k \leftarrow \phi(t_k)$ 
    if  $x_k \in L$  then
       $loss \leftarrow \frac{1}{|\beta|} \left\{ \sum_{x_k \in L \cap \beta} [J(y_k, s_k) + J(y_k, t_k)] \right\}$ 
    end if
    if  $x_k \in U$  then
       $loss \leftarrow \frac{1}{|\beta|} \left\{ \sum_{x_k \in U \cap \beta} [J(\tilde{t}_k, s_k) + a \cdot J(\tilde{s}_k, t_k)] \right\}$ 
    end if
    update  $\theta, \theta'$   $\triangleright$  update network parameters
  end for
end for
    
```

---

shown in Algorithm 1, an end-to-end process is employed to train these two models.

#### 2.4. Adaptive post-processing

The median filter is utilized for post-processing of the frame-level probabilities output by the model. Instead of determining the window size of the median filter empirically, we adopt a group of median filters with adaptive window sizes for different event categories by the following formulation based on the varying length of different event categories in real life:

$$S_{win} = duration_{ave} \cdot \beta \quad (11)$$

All the frame-level probabilities output by the network are smoothed by a group of median filters with these adaptive window sizes. After smoothed, the probabilities are converted into the 0-1 prediction with a threshold of 0.5 as described in Section 2.1. Then the operation of smoothing is repeated again on the final frame-level prediction.

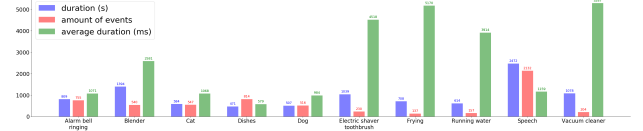


Figure 3: The total duration, number of events and average duration per event category in the synthetic training set.

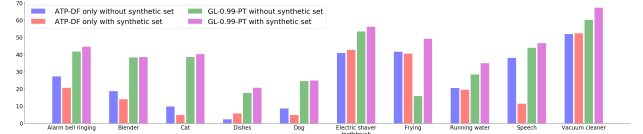


Figure 4: The class-wise  $F_1$  performance

### 3. EXPERIMENTS

#### 3.1. DCASE 2019 Task 4 Dataset

The dataset [5, 6, 7, 8, 9] of DCASE2019 task 4 is divided into 4 subsets: the weakly annotated training set (1578 clips), the unlabeled training set (14412 clips), the strongly annotated validation set (1168 clips) and the strongly annotated synthetic training set (2045 clips) [10]. We integrate the weakly annotated training set, the unlabeled training set and the strongly annotated synthetic training set (actually we only use weakly labels during training) into a training set and take the validation set as our validation set. The average duration of each event category in the synthetic set is shown in Figure 3.

#### 3.2. Feature exaction and post-processing

We produce 64 log mel-bank magnitudes which are extracted from 40 ms frames with 50% overlap ( $n_{FFT} = 2048$ ) using librosa package [11]. All the 10-second audio clips are extracted to feature vectors with 500 frames. In post-processing, we take  $\beta = \frac{1}{3}$  (see details in Section 2.4) in our experiments and the window sizes for different events are shown in Table 1.

#### 3.3. Model architecture

The model architectures of the PS-model and PT-model are described in detail in Section 2. We take  $m = 0.04$  for DF. The dimension of DF per category is shown in Table 1. The PS-model has about 2.6 times the number of trainable parameters as the PT-model (877380 / 332364). The start epoch for GL is set to 5. The PS-model with only weakly-supervised learning is named ATP-DF and the co-teaching of the PS-model and the PT-model is named GL- $\alpha$ -PT in the performance report, where  $\alpha$  is a hyper-parameter for GL discussed in Algorithm 1.

#### 3.4. Training

The Adam optimizer [12] with learning rate of 0.0018 and mini-batch of 64 10-second patches is utilized to train models. The learning rate is reduced by 20% per 10 epochs. Training early stop if there is no more improvement on clip-level macro  $F_1$  within 20 epochs. All the experiments are repeated 30 times and we report the average results. Event-based measures [13] with a 200ms collar on onsets and a 200ms / 20% of the events length collar on offsets are calculated over the entire test set.

Table 1: The dimension of the disentangled feature when  $m = 0.04$  and the window sizes of the median filters when  $\beta = \frac{1}{3}$ .

Event	DF dimension	Window size (frame)
Alarm/bell/ringing	137	17
Blender	94	42
Cat	134	17
Dishes	69	9
Dog	132	16
Electric shaver/toothbrush	76	74
Frying	34	85
Running water	160	64
Speech	30	18
Vacuum cleaner	113	87

Table 2: The performance of models from top1 and the ensemble of models.

Model	Macro $F_1$ (%)	
	Event-based	Segment-based
Top1	44.47	66.74
Ensemble (Top1-6)	45.28	<b>69.06</b>
Ensemble (Top2-6)	<b>45.43</b>	69.02

### 3.5. Results

As shown in Table 3, GL-0.99-PT (with synthetic set) achieves the best average performance on event-based macro  $F_1$ . The class-wise  $F_1$  performance per event category is shown in Figure 3. As shown in Table 2, the ensemble of the models (GL-0.99-PT) from top2 to top6 achieves the best performance, improving the performance by 21.73 percentage points from the baseline. As shown in Table 3, all the models with semi-supervised learning outperform those only with weakly-supervised learning significantly and the model with the best average performance improves the performance by 20.67 percentage points from the weakly-supervised only method. As shown in Figure 5, the performances of all the models without disentangled feature or adaptive window sizes are poorer than those which has.

#### 3.5.1. Does the synthetic training set help?

As shown in Table 3, when learning only with weakly labeled data, the synthetic training set not only does not help improve the results but also brings negative effects. But when combining weakly-supervised learning with semi-supervised learning, the synthetic training set contributes a lot so that the performance is raised by about 5-8 percentage points. We argue that the model tends to be overfitting in the synthetic training set and have difficulty in recognizing the audio clips from the real-life recording since the number of audio clips in the synthetic training set is almost 1.3 times as much as that in the weakly annotated training set. However, the large scale of unlabeled data complements this weakness and enable the synthetic training set to play a positive role during training.

#### 3.5.2. Challenge results

The model (Ensemble Top1-6) achieves an F-measure of 42.7% on the test set and won the first price in the challenge, which is 0.6 percentage point ahead of the second place and 0.7 percentage points

Table 3: The performance of models

Model	Macro $F_1$ (%)	
	Event-based	Segment-based
baseline	23.7	55.2
without the synthetic training set		
ATP-DF	25.95 ± 3.22	56.82 ± 1.34
GL-1-PT	35.19 ± 3.86	61.14 ± 3.14
GL-0.996-PT	<b>36.50 ± 3.71</b>	<b>62.03 ± 3.25</b>
GL-0.99-PT	36.21 ± 4.63	61.25 ± 2.77
GL-0.98-PT	33.78 ± 2.95	57.54 ± 3.42
with the synthetic training set		
ATP-DF	21.65 ± 2.55	57.02 ± 1.93
GL-1-PT	41.03 ± 2.98	65.58 ± 2.84
GL-0.996-PT	42.02 ± 3.29	<b>66.62 ± 1.82</b>
GL-0.99-PT	<b>42.32 ± 2.21</b>	65.78 ± 2.63
GL-0.98-PT	41.16 ± 2.42	63.89 ± 2.20

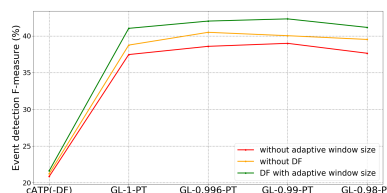


Figure 5: The performance of models without the disentangled feature (without DF) or smoothed by the median filter with a fixed window size of 27 (without adaptive window size).

ahead of the third place. We note that according to the supplementary metrics released by challenge official, our model achieves the best performance on the Youtube dataset but shows a poorer performance than the second place and the third on the Vimeo dataset. This might be because most of the audios in the dataset are from Youtube. From this point, we guess the data augmentation such as pitch shifting and time stretching might help a lot.

## 4. CONCLUSIONS

In this paper, we present a system for DCASE2019 task 4. Actually, we present a complete system for large-scale weakly labeled semi-supervised sound event detection in domestic environments. We broke the task down into 4 small sub-problems and came up with solutions for each. We release the implement to reproduce our system at [https://github.com/Kikyo-16/Sound\\_event\\_detection](https://github.com/Kikyo-16/Sound_event_detection). We employ a CNN model with an embedding-level attention module to carry out weakly-supervised learning and utilize GL to carry out semi-supervised learning. DF is employed to raise the performance of the model by reducing the interference caused by the co-occurrence of multiple event categories. In addition, adaptive post-processing is proposed to get more accurate detection boundaries. We also analyze the effect of the synthetic training set. As a result, we achieve state-of-the-art performance on the dataset of DCASE2019 task4.

## 5. ACKNOWLEDGMENT

This work is partly supported by Beijing Natural Science Foundation (4172058).

## 6. REFERENCES

- [1] R. Serizel, N. Turpault, H. Eghbal-Zadeh, and A. P. Shah, “Large-scale weakly labeled semi-supervised sound event detection in domestic environments,” in *Proceedings of the Detection and Classification of Acoustic Scenes and Events 2018 Workshop (DCASE2018)*, November 2018, pp. 19–23. [Online]. Available: <https://hal.inria.fr/hal-01850270>
- [2] L. Lin, X. Wang, H. Liu, and Y. Qian, “Specialized decision surface and disentangled feature for weakly-supervised polyphonic sound event detection,” *arXiv preprint arXiv:1905.10091*, 2019.
- [3] —, “Guided learning for the combination of weakly-supervised and semi-supervised learning,” *arXiv preprint arXiv:1906.02517*, 2019.
- [4] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” *arXiv preprint arXiv:1502.03167*, 2015.
- [5] N. Turpault, R. Serizel, A. Parag Shah, and J. Salamon, “Sound event detection in domestic environments with weakly labeled data and soundscape synthesis,” June 2019, working paper or preprint. [Online]. Available: <https://hal.inria.fr/hal-02160855>
- [6] J. F. Gemmeke, D. P. W. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, “Audio set: An ontology and human-labeled dataset for audio events,” in *Proc. IEEE ICASSP 2017*, New Orleans, LA, 2017.
- [7] E. Fonseca, J. Pons, X. Favory, F. Font, D. Bogdanov, A. Ferraro, S. Oramas, A. Porter, and X. Serra, “Freesound datasets: a platform for the creation of open audio datasets,” in *Proceedings of the 18th International Society for Music Information Retrieval Conference (ISMIR 2017)*, Suzhou, China, 2017, pp. 486–493.
- [8] F. Font, G. Roma, and X. Serra, “Freesound technical demo,” in *Proceedings of the 21st ACM international conference on Multimedia*. ACM, 2013, pp. 411–412.
- [9] G. Dekkers, S. Lauwereins, B. Thoen, M. W. Adhana, H. Brouckxon, T. van Waterschoot, B. Vanrumste, M. Verhelst, and P. Karsmakers, “The SINS database for detection of daily activities in a home environment using an acoustic sensor network,” in *Proceedings of the Detection and Classification of Acoustic Scenes and Events 2017 Workshop (DCASE2017)*, November 2017, pp. 32–36.
- [10] J. Salamon, D. MacConnell, M. Cartwright, P. Li, and J. P. Bello, “Scaper: A library for soundscape synthesis and augmentation,” in *2017 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2017, pp. 344–348.
- [11] Brian McFee, Colin Raffel, Dawen Liang, Daniel P.W. Ellis, Matt McVicar, Eric Battenberg, and Oriol Nieto, “librosa: Audio and Music Signal Analysis in Python,” in *Proceedings of the 14th Python in Science Conference*, Kathryn Huff and James Bergstra, Eds., 2015, pp. 18 – 24.
- [12] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [13] A. Mesaros, T. Heittola, and T. Virtanen, “Metrics for polyphonic sound event detection,” *Applied Sciences*, vol. 6, no. 6, p. 162, 2016. [Online]. Available: <http://www.mdpi.com/2076-3417/6/6/162>