

# Leveraging Audio-Tagging Assisted Sound Event Detection using Weakified Strong Labels and Frequency Dynamic Convolutions

Tanmay Khandelwal\*, Rohan Kumar Das\*, Andrew Koh† and Eng Siong Chng†

\*Fortemedia Singapore, Singapore †Nanyang Technological University, Singapore

Email: f20170106p@alumni.bits-pilani.ac.in, rohankd@fortemedia.com, andr0081@e.ntu.edu.sg, aseschn@ntu.edu.sg

**Abstract**—Jointly learning from a small labeled set and a larger unlabeled set is an active research topic under semi-supervised learning (SSL). In this paper, we propose a novel SSL method based on a two-stage framework for leveraging a large unlabeled in-domain set. Stage-1 of our proposed framework focuses on audio-tagging (AT), which assists the sound event detection (SED) system in Stage-2. The AT system is trained utilizing a strongly labeled set converted into weak predictions referred to as weakified set, a weakly labeled set, and an unlabeled set. This AT system then infers on the unlabeled set to generate reliable pseudo-weak labels, which are used with the strongly and weakly labeled set to train a frequency dynamic convolutional recurrent neural network-based SED system at Stage-2 in a supervised manner. Our system outperforms the baseline by 45.5% in terms of polyphonic sound detection score on the DESED real validation set.

**Index Terms**—semi-supervised learning, sound event detection, two-stage setup, pseudo-labels

## I. INTRODUCTION

Sound aids us in perceiving environmental changes and comprehending our surroundings. Humans have an in-built system for detecting and categorizing sound events in our various environments. The SED applications include audio surveillance in a variety of environments, such as smart-homes [1], [2] and cities [3], [4].

The models developed for SED require strongly labeled data to accurately predict the temporal onset and offset. The manual annotation process for generating strong labels is expensive and time-consuming, and the annotations vary greatly due to the subjective judgment of the annotators. On the other hand, annotating the entire clip with audio labels to generate weak labels is much easier. Furthermore, collecting unlabeled datasets in-domain is equally simple. To leverage this readily available unlabeled set with a small amount of labeled set, several previous works have employed semi-supervised learning (SSL) techniques. The authors of [5] used mean-teacher (MT) learning method that employs exponential moving average, whereas an unsupervised data-augmentation is used in [6], which enforces the model to be consistent with respect to the noise added using data-augmentation (DA) techniques. In [7], the authors used interpolation consistency training (ICT) and shift consistency training, whereas in [8], they self-trained to produce pseudo-labels and train on them.

To use the unlabeled set for supervised learning, the model generates pseudo-labels [9]–[11]. The pseudo-labeling process

is similar to entropy minimization [12] and helps in cases where it can recover the cluster structure among the various classes [13]. It requires a sufficient number of labeled points to effectively learn the differentiation between the clusters. The labels for pseudo-labels are determined by the confidence threshold. The clip-wise labels above the confidence threshold are used as true labels for clips [11] in the typical supervised loss function. The model is then trained using labeled and unlabeled sets simultaneously. In addition to SSL techniques, past works have employed various DA techniques like SpecAugment [14], time-shift [9], pitch-shift [15], and mixup [16] to increase diversity and reduce overfitting.

The detection and classification of acoustic scenes and events (DCASE) 2022 Task 4 focuses on SED-based SSL to utilize labeled and unlabeled data. We make the following contributions in this work to effectively exploit the unlabeled in-domain set provided in DCASE 2022 Task 4 by generating pseudo-labels:

- Proposal of a two-stage framework [17] that performs audio-tagging (AT) at the first stage to estimate reliable pseudo-labels on unlabeled data used to train the SED system at the second stage in a supervised manner.
- A novel weak training strategy to create weakified labels, where the strong labels are converted to weak predictions. The objective is to supply more weak labels for Stage-1 system training to lessen the model’s inclination to predict inactive frames [18] when trained with strong labels.
- Utilize pre-trained audio neural networks (PANNs) for Stage-1 to further improve the reliability of the pseudo-weak labels used to train the Stage-2 system.

We used several DA techniques, pooling functions, and adaptive post-processing to improve the robustness of the developed systems, evaluate the proposed method, and make fair comparisons with other state-of-the-art methods.

## II. SOUND EVENT DETECTION SYSTEM

In this section, we briefly review the baseline system and the proposed, two-stage framework for sound event detection.

### A. Baseline

The baseline [19] architecture is a combination of convolutional neural network (CNN) and recurrent neural network (RNN) called convolutional recurrent neural network (CRNN),

as depicted in Fig. 1 (a). The CNN part is made up of 7-blocks, each with 16, 32, 64, 128, 128, 128, and 128 filters, respectively. It has a kernel size of  $3 \times 3$  and an average-pooling of [2, 2], [2, 2], [2, 1], [2, 1], [2, 1], [2, 1], [2, 1] per layer. The RNN is composed of two layers of 128 bidirectional gated recurrent units (Bi-GRU) [20]. The RNN block is followed by an attention pooling layer, which is a multiplication of a linear layer with softmax activations and a linear layer with sigmoid activations. The baseline employs the MT [5] strategy, which is a hybrid of two models: the student model and the teacher model (both having the same architecture). The student model is the final model used for inference, whereas the teacher model is designed to help the student model during training. Its weights are an exponential moving average of the student model’s weights.

### B. Proposed two-stage framework

The objective of any SSL algorithm is to utilize labeled and unlabeled data to learn the underlying structure of the dataset effectively. The small amount of labeled data helps the model learn discrete or non-overlapping clusters for different labels. The cluster assumption [7] states that close points have the same class and points in different classes are more widely separated, therefore true decision boundaries flow through low-density input space. As training progresses, these clusters improve their cluster boundary, improving model predictions on the unlabeled in-domain set. To mitigate the problem of a small labeled set and to utilize the unlabeled set by learning the discrete clusters for each class, we propose a two-stage framework [17], shown in Fig. 2. Stage-1 utilizes the proposed weak training method to focus on AT, and Stage-2 then utilizes the reliable pseudo-labels generated from Stage-1 to have an improved SED performance. Furthermore, each stage makes use of MT adopted from the baseline. In addition to MT, we use another method used for SSL, called ICT [21], in both stages of the two-stage framework. The ICT substitutes all input samples with interpolated samples, helping the model to improve the generalization ability. A detailed description of the models used in each stage is given in the following subsections.

1) *Stage-1*: In order to have an effective AT in Stage-1, [9] showed the importance of deeper neural network models compared to the baseline CRNN. As feature extractor, we used CNN-14-based PANNs [22] with 118M parameters for pre-trained embeddings. The parameters of the PANNs-based embeddings are unfrozen and trained. The 14-layer CNN feature extractor consists of 6 convolutional blocks. Each convolutional block consists of 2 convolutional layers with a kernel size of  $3 \times 3$ . In addition, each convolutional layer is followed by batch normalization and rectified linear unit [23] non-linearity to stabilize the training. Average pooling [24] of  $2 \times 2$  is applied to each convolutional block for down-sampling. The feature extractor is followed by 2-layers of Bi-GRU with 1024 hidden units. For frame-level predictions, the RNN output is multiplied by a dense layer with sigmoid activation, and for clip-level predictions, the linear layer is

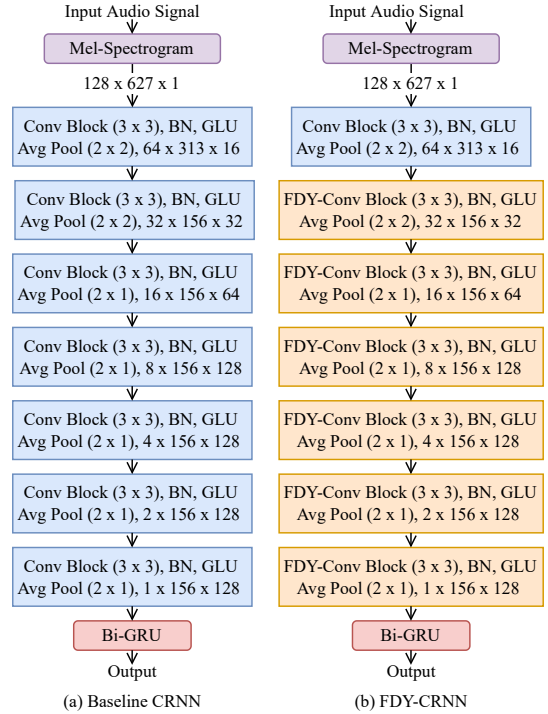


Fig. 1. Architecture of (a) CRNN (Baseline) (b) FDY-CRNN.

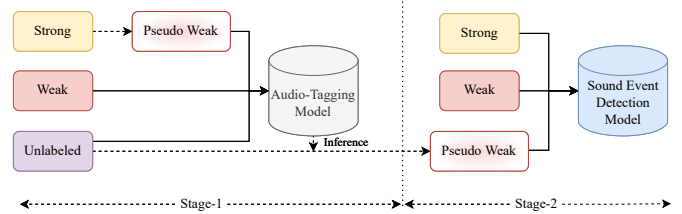


Fig. 2. Proposed two-stage learning setup, with Stage-1 focusing on AT and Stage-2 focusing on SED.

multiplied by a dense layer with softmax activation. Based on previous work [18] that uses only the weakly labeled data, we suggest a weak training strategy to improve Stage-1 AT systems. We converted the strongly labeled set into a weakly labeled set by removing the onset and offset and keeping the event labels, which we refer to as *weakified labels*. Then we trained the AT system using the weakified labels, weakly labeled set, and unlabeled set as illustrated in Fig. 2.

2) *Stage-2*: In this work, we used the AT (Stage-1) based system to make predictions on the unlabeled set to use them as pseudo-weak labels in Stage-2 training, as shown in Fig. 2. We believe this way, we can generate reliable pseudo-labels, which can help the SED model at Stage-2. The baseline CRNN’s standard 2D convolutional block enforces translation equivariance on sound events along both the time and frequency axes, despite the fact that frequency is not a shift-invariant dimension. To focus on frequency-dependent patterns and to further improve the SED performance, we employed frequency dynamic (FDY)-convolutions proposed in [25] as it applies

frequency adaptive kernels to enforce frequency dependency on 2D convolution. We replaced the baseline’s standard 2D convolutional blocks with FDY-convolutional blocks, which have the same number of layers and feature maps as that of the baseline, as illustrated in Fig. 1 (b). Then it was trained on a pseudo-weakly labeled set, in addition to the strongly labeled set and the weakly labeled set, in a supervised manner.

### III. ADDITIONAL METHODS

#### A. Pooling function

Motivated from a prior work [26], we used exponential softmax to replace the attention pooling used in the baseline. The exponential softmax function assigns a weight of  $\exp(y_i)$  to the frame-level probability  $y_i$  as given below:

$$y = \frac{\sum_i y_i \exp(y_i)}{\sum_i \exp(y_i)} \quad (1)$$

where  $y_i$  is the predicted probability of an event occurring in the  $i^{\text{th}}$  frame. This implies that, with a higher prediction probability, the higher the exponential weight is assigned to the frame-level probability. Hence, it is better under the stringent evaluation criteria for the correctness of the category.

#### B. Asymmetric focal loss (AFL)

AFL [27] function is used to control the training weight depending on the ease and difficulty of the model training. The AFL function for each  $k^{\text{th}}$  data point with target sound event as  $y_k$  and predicted sound event as  $p_k$  is given below:

$$l_{AFL}(p, y) = \sum_{n=1}^K [(1-p_k)^\gamma y_k \ln p_k + (p_k)^\zeta (1-y_k) \ln(1-p_k)] \quad (2)$$

where the parameters  $\gamma$  and  $\zeta$  are the weighing hyperparameters given as the input to the function that controls the weight of active and inactivate frames.

#### C. Data-augmentation (DA)

We used several DA techniques during the training in both stages, such as time-masking [14], frame-shifting [18], mixup [16], addition of Gaussian noise and filter augmentation [18]. Time-masking masks the sequential time frames (which means replacing the elements by zeros or other values), whereas frame-shifting shifts the features and labels along the time axis. Again, mixup randomly mixes selected samples with a mixing parameter, helping in linear interpolation to improve the robustness of the model. In addition, filter augmentation, which uses varying weights on random frequency regions, has been shown to significantly improve SED performance.

#### D. Adaptive post-processing

We used adaptive post-processing [28] in all trials, where the median filter window sizes ( $Win$ ) are different for each event category  $c$  based on the real-life event lengths as shown below:

$$Win_c = duration_c \times \beta_c \quad (3)$$

For several event categories with high duration variance, we used  $duration_c$  as the median duration. Here, we used  $\beta_c = \frac{1}{3}$  and then slightly adjusted the window sizes on the validation set.

## IV. EXPERIMENTAL SETUP

The next subsections outline our experimental setup to demonstrate the efficacy of the proposed methods.

#### A. Dataset

The DCASE 2022 Task 4 dataset used in this work is composed of 10 seconds audio clips to simulate a domestic environment. The development training set is divided into 3 major subsets:

- 1,578 real recordings with weak annotations.
- 14,412 real recordings, unlabeled in-domain training set
- 10,000 synthetic recordings with strong annotations [29].
- An additional subset from the recently released strongly labeled AudioSet [30] subset of 3,470 real recordings with strong annotations is released as external data.

The development validation set has 1,168 real recordings with strong annotations, and the public evaluation (“YouTube”) set has 692 YouTube clips.

#### B. Pre-processing

The audio clips are re-sampled at 16 kHz to a mono channel. Then, log-mel spectrograms are produced using mel-filters in the frequency domain from 0 to 8 kHz with a window size of 2048 samples and a hop size of 256 samples. Stage-1 employed 64 mel-filters, while Stage-2 used 128 mel-filters. The clips with a duration of less than 10 seconds are padded with silence.

#### C. Training process

The batch size for all the experiments is 48 (1/4 strong set, 1/4 weak set, 1/2 unlabeled set). We employed Adam optimizer [31] with a learning rate of 0.001 and an exponential warmup for the first 50 epochs with no early stopping.

#### D. Evaluation metrics

We used polyphonic sound event detection scores (PSDS) [32] as a performance metric in our studies. The PSDS is more resistant to labeling subjectivity, allowing for ground truth interpretation and temporal structure detection. The single PSDS is computed using polyphonic receiver operating characteristic curves, allowing comparison independent of the operating point. Additionally, it can be adapted for various applications to ensure the appropriate user experience. As a result, it overcomes the limitations of traditional collars-based event F-scores. Using hyperparameters values adopted from the DCASE 2022 Task 4 for the Detection Tolerance Criterion ( $\rho_{DTC}$ ) and Ground Truth intersection Criterion ( $\rho_{GTC}$ ) mentioned in Table I, we compute the PSDS on two scenarios that stress distinct system features. The system must react fast to event detection in Scenario-1, hence it focuses on sound event temporal localization. Scenario-2 focuses less on reaction time and more on class confusion.

TABLE I  
PSDS HYPERPARAMETERS FOR EACH EVALUATION SCENARIO.

Scenarios	$\rho_{DTC}$	$\rho_{GTC}$
Scenario-1	0.7	0.7
Scenario-2	0.1	0.1

TABLE II  
DESCRIPTION OF THE TWO-STAGE SYSTEM DEVELOPED IN THIS WORK.

Stage	DA	Description
1	time-masking, frame-shifting, mixup, and Gaussian noise addition	We used the architecture given in Section V-A, trained on weak, unlabeled, and weakified set using exponential softmax function during inference to get the best results.
2	time-masking, frame-shifting, mixup and filter augmentation	Stage-1 inferred on the unlabeled set, while Stage-2 used AFL function with $\gamma=0.625$ and $\zeta=1$ to train the architecture specified in Section V-B on weak, pseudo-weak, and strong sets.

### E. Developed system

The models, DA methods, and experimental settings of our two-stage system are given in Table II. In our two-stage study, Stage-1 uses PANNs while Stage-2 uses FDY-CRNN.

## V. RESULTS AND ANALYSIS

In this section, we report the studies of the proposed two-stage framework in a stage-wise manner, with ablation studies.

### A. Stage-1 comparison

We are first interested in assessing the contribution of each component to our system at Stage-1. Table III shows the SED performance of Stage-1 trained on a real strong set, synthetic strong set, weak set, unlabeled set, and using CNN-14-based PANNs as the pre-trained embeddings. Experimental results show that pre-trained models as feature extractors trained on larger datasets exceed the DCASE 2022 Task 4 organizers’ baseline (Baseline), which uses an external dataset. We also observe that our weak training with the PANNs method significantly improves PSDS2 from 0.552 to 0.831 compared to the baseline and drastically decreases PSDS1 from 0.351 to 0.057. PSDS2 increases due to low tolerance in  $\rho_{DTC}$  and  $\rho_{GTC}$  [32], as seen in Table I. As per the parameters for PSDS2, the tolerance value is 0.1, thus the prediction is regarded as true positive even when there is at least one ground truth greater than 1 second out of the 10 seconds clip [18]. Thus, having a higher PSDS2 is equivalent to having a better AT system.

This relation is extended to train Stage-2 using weak pseudo-labels from PANNs-based Stage-1 with a higher PSDS2. Further, the decrease in PSDS1 can be attributed to it specifically focusing on temporal localization with a tolerance value of 0.7. Using a pre-trained model also sped up training because the model converged faster with optimized weights. The results show that DA approaches mentioned in Table II and adaptive post-processing improve performance slightly. We also demonstrate the performance of our Stage-1 on the public evaluation set later in Table VI.

### B. Stage-2 comparison

To assess the two-stage framework’s importance, we constructed the baseline (CRNN) system from the organizers in a

TABLE III  
PERFORMANCE OF THE BASELINE AND OUR STAGE-1 SYSTEM ON THE REAL VALIDATION SET OF DCASE 2022 TASK 4.

Model	Method	PSDS1	PSDS2
CRNN	DCASE 2022 Task 4 Baseline	0.351	0.552
PANNs	CRNN replaced by CNN-14 PANNs	0.450	0.716
PANNs	+ Weak Training	0.057	0.831
PANNs	+ ICT	0.067	0.834
PANNs	+ DA + Post-processing	0.075	<b>0.840</b>

TABLE IV  
PERFORMANCE OF THE BASELINE AND OUR STAGE-2 SYSTEM ON THE REAL VALIDATION SET OF DCASE 2022 TASK 4.

Stage-1	Stage-2	PSDS1	PSDS2
CRNN	CRNN	0.378	0.578
PANNs	CRNN	0.437	0.681
PANNs	FDY-CRNN	0.450	0.701
PANNs	FDY-CRNN + DA	0.468	0.714
PANNs	FDY-CRNN + DA + Post-processing	0.470	0.718
PANNs	FDY-CRNN + DA + Post-processing + AFL	<b>0.472</b>	0.721

two-stage setup. Table IV demonstrates a 5.8% improvement in total PSDS (PSDS1 + PSDS2) over its result without a two-stage setup in Table III. We then used the best Stage-1 CNN-14-based PANNs model (PSDS2 = 0.840) to infer on the unlabeled data to create pseudo-weak labels for Stage-2. Using CRNN in Stage-2, resulted in a PSDS1 of 0.437 and PSDS2 of 0.681. Replacing with FDY-convolutions in Stage-2 improved the SED performance, resulting in a PSDS1 of 0.450, and with additional methods resulted in a PSDS1 of 0.472, a 45.5% improvement in overall PSDS (PSDS1 + PSDS2) for the two-stage setup. Table V shows the comparison of our system with other single systems (without ensembling) on the real validation set. The same Stage-2’s PSDS1 was 0.479 and PSDS2 was 0.733 on the public evaluation set, as shown in Table VI.

TABLE V  
PERFORMANCE COMPARISON OF SINGLE SYSTEMS ON THE REAL VALIDATION SET OF DCASE 2022 TASK 4 WITH OTHER TEAMS

Team	#Parameters	PSDS1	PSDS2	PSDS1 + PSDS2
Ebbers-UPB-task4	779M	0.505	0.807	1.312
<b>Ours</b>	2.8M	0.472	0.721	<b>1.193</b>
Zhang-UCAS-task4	11M	0.459	0.672	1.131
Kim-GIST-task4	1M	0.455	0.670	1.125
Dinkel-XiaoRice-task4	37M	0.425	0.644	1.069

TABLE VI  
PERFORMANCE ON PUBLIC EVALUATION SET OF DCASE 2022 TASK 4.

Stage-1	Stage-2	PSDS1	PSDS2
Baseline	-	0.387	0.592
PANNs	-	0.087	<b>0.801</b>
PANNs	FDY-CRNN	<b>0.479</b>	0.733

## VI. CONCLUSION

In this work, we proposed an AT-assisted sound event detection system using a two-stage framework. We introduced a weak training method to derive weakified labels from strong labels for AT system at Stage-1 and used FDY convolutions in the baseline to focus on the frequency-dependent patterns. Additionally, CNN-14-based pre-trained audio neural networks were used as pre-trained embeddings in Stage-1 to generate reliable pseudo-weak labels to utilize in Stage-2. The studies on DCASE 2022 Task 4 validation set and public evaluation set proved the importance of the proposed two-stage setup and the usage of a weak training strategy. We outperform the DCASE 2022 baseline by 45.5% on the real validation set in both aspects of the PSDS metric.

## REFERENCES

- [1] Juan Pablo Bello, Claudio Silva, Oded Nov, R. Luke Dubois, Anish Arora, Justin Salamon, Charles Mydlarz, and Harish Doraiswamy, “SONYC: A system for monitoring, analyzing, and mitigating urban noise pollution,” *Communications of the ACM*, vol. 62, no. 2, pp. 68–77, 2019.
- [2] Tanmay Khandelwal, Rohan Kumar Das, and Eng Siong Chng, “Is your baby fine at home? Baby cry sound detection in domestic environments,” *Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pp. 275–280, 2022.
- [3] Juan Pablo Bello, Charlie Mydlarz, and Justin Salamon, “Sound analysis in smart cities,” *Springer International Publishing*, pp. 373–397, 2018.
- [4] Christian Debes, Andreas Merentitis, Sergey Sukhanov, Maria Niessen, Nikolaos Frangiadakis, and Alexander Bauer, “Monitoring activities of daily living in smart homes: Understanding human behavior,” *IEEE Signal Processing Magazine*, vol. 33, no. 2, pp. 81–94, 2016.
- [5] Antti Tarvainen and Harri Valpola, “Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results,” *International Conference on Neural Information Processing Systems*, pp. 1195–1204, 2017.
- [6] Heinrich Dinkel, Xinyu Cai, and Zhiyong Yan, “The smallrice submission to the DCASE 2021 task 4 challenge: A lightweight approach for semi-supervised sound event detection with unsupervised data augmentation,” *Detection and Classification of Acoustic Scenes and Events (DCASE) Challenge*, 2021.
- [7] Xu Zheng, Han Chen, and Yan Song, “Zheng USTC team’s submission for DCASE 2021 task 4 - semi-supervised sound event detection,” *Detection and Classification of Acoustic Scenes and Events (DCASE) Challenge*, 2021.
- [8] Janek Ebberts and Reinhold Haeb-Umbach, “Self-trained audio tagging and sound event detection in domestic environments,” *Detection and Classification of Acoustic Scenes and Events (DCASE) Workshop*, pp. 226–230, 2021.
- [9] Chih-Yuan Koh, You-Siang Chen, Yi-Wen Liu, and Mingsian R. Bai, “Sound event detection by consistency training and pseudo-labeling with feature-pyramid convolutional recurrent neural networks,” *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 376–380, 2021.
- [10] Chungo Park, Donghyeon Kim, and Hanseok Ko, “Sound event detection by pseudo-labeling in weakly labeled dataset,” *Sensors*, vol. 21, no. 24, 2021.
- [11] Nam Kyun Kim and Hong Kook Kim, “Self-training with noisy student model and semi-supervised loss function for DCASE 2021 challenge task 4,” *Detection and Classification of Acoustic Scenes and Events (DCASE) Challenge*, 2021.
- [12] Yves Grandvalet and Yoshua Bengio, “Semi-supervised learning by entropy minimization,” *International Conference on Neural Information Processing Systems*, 2004.
- [13] Dong-Hyun Lee, “Pseudo-label : The simple and efficient semi-supervised learning method for deep neural networks,” *International Conference on Machine Learning (ICML)*, 2013.
- [14] Daniel Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le, “SpecAugment: A simple data augmentation method for automatic speech recognition,” *Interspeech*, pp. 2613–2617, 2019.
- [15] Brian McFee, Eric J. Humphrey, and Juan Pablo Bello, “A software framework for musical data augmentation,” *International Society for Music Information Retrieval (ISMIR)*, pp. 248–254, 2015.
- [16] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz, “Mixup: Beyond empirical risk minimization,” *International Conference on Learning Representations (ICLR)*, 2018.
- [17] Tanmay Khandelwal, Rohan Kumar Das, Andrew Koh, and Eng Siong Chng, “FMSG-NTU submission for DCASE 2022 Task 4 on sound event detection in domestic environments,” *Detection and Classification of Acoustic Scenes and Events (DCASE) Challenge*, 2022.
- [18] Hyeonuk Nam, Byeong-Yun Ko, Gyeong-Tae Lee, Seong-Hu Kim, Won-Ho Jung, Sang-Min Choi, and Yong-Hwa Park, “Heavily augmented sound event detection utilizing weak predictions,” *Detection and Classification of Acoustic Scenes and Events (DCASE) Challenge*, 2021.
- [19] Lionel Delphin-Poulat and Cyril Plapous, “Mean teacher with data augmentation for DCASE 2019 task 4 technical report,” *Detection and Classification of Acoustic Scenes and Events (DCASE) Challenge*, 2019.
- [20] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio, “Learning phrase representations using RNN encoder–decoder for statistical machine translation,” *Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1724–1734, 2014.
- [21] Vikas Verma, Alex Lamb, Juho Kannala, Yoshua Bengio, and David Lopez-Paz, “Interpolation consistency training for semi-supervised learning,” *International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 3635–3641, 2019.
- [22] Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D. Plumbley, “PANNs: Large-scale pretrained audio neural networks for audio pattern recognition,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 2880–2894, 2020.
- [23] Abien Fred Agarap, “Deep learning using rectified linear units (ReLU),” *CoRR*, vol. abs/1803.08375, 2018.
- [24] Chen-Yu Lee, Patrick Gallagher, and Zhuowen Tu, “Generalizing pooling functions in CNNs: Mixed, gated, and tree,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 4, pp. 863–875, 2017.
- [25] Hyeonuk Nam, Seong-Hu Kim, Byeong-Yun Ko, and Yong-Hwa Park, “Frequency dynamic convolution: Frequency-adaptive pattern recognition for sound event detection,” *Interspeech*, pp. 2763–2767, 2022.
- [26] Yun Wang, Juncheng Billy Li, and Florian Metzke, “A comparison of five multiple instance learning pooling functions for sound event detection with weak labeling,” *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 31–35, 2019.
- [27] Keisuke Imoto, Sakiko Mishima, Yumi Arai, and Reishi Kondo, “Impact of sound duration and inactive frames on sound event detection performance,” *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 860–864, 2021.
- [28] Koichi Miyazaki, Tatsuya Komatsu, Tomoki Hayashi, Shinji Watanabe, Tomoki Toda, and Kazuya Takeda, “Convolution-augmented transformer for semi-supervised sound event detection,” *Detection and Classification of Acoustic Scenes and Events (DCASE) Challenge*, 2020.
- [29] Francesca Ronchini, Romain Serizel, Nicolas Turpault, and Samuele Cornell, “The impact of non-target events in synthetic soundscapes for sound event detection,” *Detection and Classification of Acoustic Scenes and Events (DCASE) Workshop*, pp. 115–119, 2021.
- [30] Jort Florent Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, Channing Moore, Manoj Plakal, and Marvin Ritter, “AudioSet: An ontology and human-labeled dataset for audio events,” *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 776–780, 2017.
- [31] Diederik P. Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *International Conference on Learning Representations (ICLR)*, 2015.
- [32] Cagdas Bilen, Giacomo Ferroni, Francesco Tuveri, Juan Azcarreta, and Sacha Krstulovic, “A framework for the robust evaluation of sound event detection,” *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 61–65, 2020.