

# IMPROVING WEAKLY SUPERVISED SOUND EVENT DETECTION WITH CAUSAL INTERVENTION

Yifei Xin<sup>1</sup>, Dongchao Yang<sup>1</sup>, Fan Cui<sup>2</sup>, Yujun Wang<sup>2</sup>, Yuexian Zou<sup>1,\*</sup>

<sup>1</sup>School of ECE, Peking University, Shenzhen, China

<sup>2</sup>Xiaomi Corporation, Beijing, China

## ABSTRACT

Existing weakly supervised sound event detection (WSSSED) work has not explored both types of co-occurrences simultaneously, i.e., some sound events often co-occur, and their occurrences are usually accompanied by specific background sounds, so they would be inevitably entangled, causing misclassification and biased localization results with only clip-level supervision. To tackle this issue, we first establish a structural causal model (SCM) to reveal that the context is the main cause of co-occurrence confounders that mislead the model to learn spurious correlations between frames and clip-level labels. Based on the causal analysis, we propose a causal intervention (CI) method for WSSSED to remove the negative impact of co-occurrence confounders by iteratively accumulating every possible context of each class and then re-projecting the contexts to the frame-level features for making the event boundary clearer. Experiments show that our method effectively improves the performance on multiple datasets and can generalize to various baseline models.

**Index Terms**— Causal intervention, weakly supervised sound event detection, structural causal model

## 1. INTRODUCTION

Sound event detection (SED) involves two subtasks: one is to recognize the types of sound events in an audio clip (audio tagging), and the other is to pinpoint their onset and offset times (localization). Since frame-level labels are costly to collect, weakly supervised sound event detection (WSSSED) [1, 2] has gained an increasing research interest, which has only access to weak clip-level labels in the training stage, yet requires to perform the frame-level prediction of onset and offset times during evaluation.

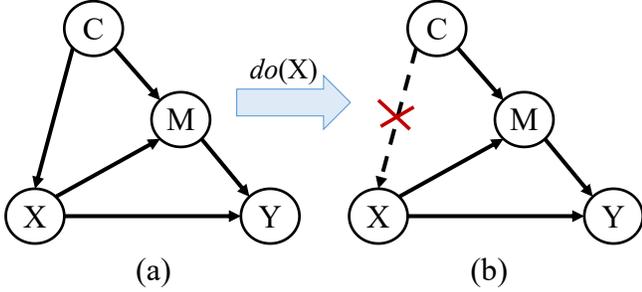
However, one challenging problem of WSSSED is that some sound events often co-occur in an audio clip (e.g., the two classes “train” and “train horn” in DCASE2017 task4

dataset [3]). As a result, it is difficult to distinguish those frequently co-occurring sound events in an audio recording, since the model will inevitably relate the sound event class of “train” with that of “train horn”, which interferes with the recognition and detection of each other. Some approaches have been proposed to address this issue. In [4], graph Laplacian regularization was introduced to model the co-occurrence of sound events for strong labeled SED. In [5], Lin et al. proposed a disentangled feature, which re-models the high-level feature space so that the feature subspace can be different for each sound event. However, the co-occurrence issue is not just between sound events as sound events also usually co-occur with specific background sounds. Thus, the sound events and background sounds would be inevitably entangled, causing the model to falsely generate ambiguous frame-level localization results with only clip-level supervision. In this work, we target the co-occurrence issue from the two aspects mentioned above for WSSSED, which we call “entangled context”. The “entangled context” would lead to misclassification and biased localization results, including the wrongly confusing co-occurring sound events and entangled background sounds. Therefore, we argue that resolving the “entangled context” issue is essential for WSSSED.

In this paper, we attempt to address this issue with causal intervention (CI) method, called CI-WSSSED. CI-WSSSED attributes the “entangled context” to the frequently co-occurring sound events and specific background sounds that mislead the model to learn spurious correlations between frames and clip-level labels. To find those frames which truly contribute to the clip-level labels in an audio clip, we first establish a structural causal model (SCM) [6] to clarify the causal relation among frame-level features, contexts, and clip-level labels, revealing that the context is the main cause of co-occurrence confounders. Ideally, if we could collect enough audio clips covering all the combinations of different sound event co-occurrences under various background sounds in a balanced distribution, we can distinguish any sound event from them easily. However, it is labor-intensive or even impossible to collect such a huge dataset for each sound event class. To this end, we employ causal intervention to intervene the input to be under any possible context in an approximate way. Based on the causal analysis, CI-WSSSED is then designed to operate

This paper was partially supported by NSFC (No: 62176008) and Shenzhen Science & Technology Research Program (No:GXWD20201231165807007-20200814115301001). Special acknowledgements are given to Xiaomi for its support.

\* Yuexian Zou is the corresponding author.



**Fig. 1.** (a) The structural causal model (SCM) for WSSSED. (b) The intervened SCM based on backdoor adjustment for WSSSED.

in an iterative procedure, which is achieved by accumulating the contextual information for each class and then employing it as attention to enhance the frame-level representation for making the sound event boundary clearer.

In summary, the contributions of this paper are as follows:

- 1) Our work is the first to concern and reveal the “entangled context” issue of WSSSED from both aspects of entangled co-occurring sound events and background sounds.
- 2) We are the first to introduce causal intervention into WSSSED for the “entangled context” issue and design a new network structure, called CI-WSSSED, to embed the causal intervention into the WSSSED pipeline with an end-to-end scheme.
- 3) Experiments show that our CI-WSSSED yields significant performance gains on WSSSED datasets and can generalize to various baselines.

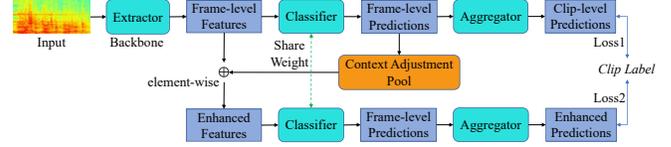
## 2. CAUSAL INTERVENTION

The goal of causal learning is to enable the model to pursue causal effects: it can eliminate the spurious bias and disentangle the desired model effects by pursuing the true causal effect [7]. Nowadays, the structural causal model (SCM) is commonly utilized in causal learning scenarios. SCM employs a graphical formalism in which nodes are represented as random variables and directed edges represent the direct causal relationship between these variables. As shown in Fig. 1(a), the conditional distribution  $P(Y|X)$  expresses the likelihood of  $Y$  given  $X$ , where  $Y$  is not only caused by  $X$  via  $X \rightarrow Y$ , but also  $C$  via the correlation  $C \rightarrow X \rightarrow Y$ .

To find the causal effect of variable  $X$  on variable  $Y$ , do-calculus is introduced [8]. In detail, causal intervention [9] fixes the target variable  $X$  to a constant  $x$ , denoted as  $do(X = x)$ , rendering it independent of its causes, so the causal effect of  $X$  on  $Y$  is formulated as:

$$P(Y|do(X = x)), \quad (1)$$

where the do-notation denotes intervening to set variable  $X$  to the value  $x$ , thereby removing all incoming arrows to the variable  $X$ , as shown in Fig. 1(b). In this way, we can lever-



**Fig. 2.** Overview of our proposed CI-WSSSED approach.

age  $P(Y|do(X))$  to pursue the true causality and remove the negative effect of confounders.

A straightforward way to intervene  $X$  is conducting a randomized controlled trial [10] with an ideal huge dataset, which contains audio clips of all kinds of sound event co-occurrences under various background sounds. Therefore, the spurious correlation  $C \rightarrow X$  is cutoff and then  $P(Y|X) = P(Y|do(X))$ . Since this kind of intervention is impossible due to the huge cost of collecting such a large dataset for each sound event, we apply the backdoor adjustment [11] to approximate it. In the following section, we will detail how we leverage this solution to WSSSED in our approach.

## 3. PROPOSED METHOD

### 3.1. Structural Causal Model for WSSSED

In this part, we will explain why the entangled context hinders the sound event classification and localization performance in WSSSED. We formulate the causalities among frame-level features  $X$ , contexts  $C$ , frame-level predictions  $M$ , and predicted clip-level labels  $Y$ , with a structural causal model (SCM) illustrated in Fig. 1(a). The direct links represent the causalities between the two nodes: cause  $\rightarrow$  effect.

$C \rightarrow X$ : We denote  $C$  as the prior context knowledge. This link represents that the extractor produces frame-level features  $X$  under the effect of contexts  $C$ . Although the contextual information improves the association between the frame-level features  $X$  and predicted labels  $Y$  via  $P(Y|X)$ ,  $P(Y|X)$  mistakenly associates non-causal but positively correlated frames to labels. There are many possible sources of the contextual bias, which may be from the training process (e.g., batch normalization) and biased datasets. For example, the sound events of “train horn” and “train” often occur in co-occurrence with each other, so the model would be confounded to establish a spurious correction between the two sound event classes, and this causal link will exacerbate the existing bias.

$C \rightarrow M \leftarrow X$ : We term the mediator  $M$  as the  $X$ -specific context, which is directly from  $X$  but essentially inherited from  $C$ . Specifically, the context is a combination of various other sound events, for instance, when multiple sound events (e.g., “children playing”, “street music”, and “dog bark”) occur in an audio clip, the “children playing” can be seen as the label with its context including “street music” and “dog bark”, and the same holds true when the lead is “street music” or “dog bark”.

**Table 1.** Performance comparison of CI-WSSSED and baseline models on the DCASE2017 task4 validation and evaluation set.

Method	Validation Set			Evaluation Set				
	AT-mAP	AT-F1	SED-mAP	AT-mAP	AT-F1	SED-mAP	Seg-F1	Event-F1
Winner SED [12]	-	-	-	-	0.526	-	0.555	-
CDur [13]	-	-	-	-	0.553	-	0.508	0.152
CNN-biGRU [1]	0.650	0.555	0.456	0.650	0.632	0.444	0.564	-
CNN-Transformer [1]	0.653	0.557	0.437	0.656	0.629	0.454	0.556	-
HTSAT [14]	0.661	0.560	0.524	0.668	0.636	0.535	0.587	0.178
CDur-CI	-	-	-	-	<b>0.561</b>	-	<b>0.511</b>	<b>0.164</b>
CNN-biGRU-CI	<b>0.661</b>	<b>0.566</b>	<b>0.462</b>	<b>0.662</b>	<b>0.641</b>	<b>0.453</b>	<b>0.570</b>	-
CNN-Transformer-CI	<b>0.662</b>	<b>0.568</b>	<b>0.449</b>	<b>0.666</b>	<b>0.637</b>	<b>0.461</b>	<b>0.561</b>	-
HTSAT-CI	<b>0.672</b>	<b>0.572</b>	<b>0.531</b>	<b>0.678</b>	<b>0.644</b>	<b>0.544</b>	<b>0.592</b>	<b>0.191</b>

$X \rightarrow Y \leftarrow M$ : These links denote that the sound event itself and its context together affect the final prediction. However, a general  $C$  cannot directly influence the predicted clip-level labels  $Y$ . Thus, in addition to the direct effect  $X \rightarrow Y$ ,  $Y$  is also the effect of the  $X$ -specific context  $M$ , which contains the timestamp information of the sound event and its context.

Considering the impact of contexts  $C$  on frame-level features  $X$ , we will cut off the link from  $C$  to  $X$ . Next, we will introduce a causal intervention method to mitigate the negative impact of contextual bias.

### 3.2. Causal Intervention via Backdoor Adjustment

As shown in Fig. 1(b), we propose to use  $P(Y|do(X))$  based on the backdoor adjustment [9, 15] to remove the context confounder and pursue the true causality between  $X$  and  $Y$ . The key idea is to cut off the link  $C \rightarrow X$ , and stratify  $C$  into pieces  $C = \{c_1, c_2, \dots, c_k\}$ , where  $c_i$  denotes the  $i^{th}$  class context. Formally, we have

$$P(Y|do(X)) = \sum_i^k P(Y|X = x, M = f(x, c_i))P(c_i), \quad (2)$$

where  $f(x, c_i)$  represents that  $M$  is formed by the combination of  $X$  and  $C$ , and  $k$  is the number of sound event classes. As  $C$  is no longer correlated with  $X$ , the causal intervention guarantees  $X$  to have an equal chance of incorporating every context  $c_i$  into  $Y$ 's prediction, based on the proportion of each  $c_i$  in the whole. However, the cost of the network forward propagation for all the  $k$  classes is expensive. Thanks to the Normalized Weighted Geometric Mean [16], we can optimize Eq. (2) to approximate the above expectation by moving the outer sum  $\sum_i^n P(c_i)$  into the feature level

$$P(Y|do(X)) \approx P(Y|X = x, M = \sum_i^k f(x, c_i)P(c_i)). \quad (3)$$

Thus, we only need to feed-forward the network once instead of  $k$  times. To simplify the formula, we assume roughly the

same number of samples for each class in the dataset, so  $P(c_i)$  is set as the uniform  $1/k$ . After simplifying Eq. (3), we have

$$P(Y|do(X)) \approx P(Y|x \oplus \frac{1}{k} \sum_i^k f(x, c_i)), \quad (4)$$

where  $\oplus$  denotes projection. Therefore, the ‘‘entangled context’’ issue has been transferred to the calculation of  $\sum_i^k f(x, c_i)$ . Next, we will introduce a context adjustment pool to represent  $\sum_i^k f(x, c_i)$ .

### 3.3. Network Structure

In this part, we implement causal intervention for WSSSED with a novel network structure, called CI-WSSSED. Fig. 2 illustrates the overview of our CI-WSSSED. First, the feature extractor (e.g., the CNN, RNN or Transformer-based backbone) takes the mel-spectrogram as input and produces high-level features  $X \in \mathbb{R}^{c \times n}$ , where  $c$  is the number of channels and  $n$  is the number of audio frames. Then, the frame-level features  $X$  are fed into the classifier with a fully connected layer followed by an aggregator (e.g., an average pooling function) to produce clip-level prediction results.

We maintain a context adjustment pool  $Q \in \mathbb{R}^{k \times n}$  for each sound event class during the training stage, which is the core of our CI-WSSSED. According to Eq. (4),  $Q$  is designed to continuously store the contextual information of each occurring sound event, and then re-project the accumulated contexts onto the frame-level features  $X$  generated by the backbone to produce enhanced features  $X^e \in \mathbb{R}^{c \times n}$ . In detail, the context adjustment pool is updated by fusing the contexts of each occurring sound event in frame-level predictions  $M = f(x, c_i)$ , which is followed by a batch normalization:

$$Q_j = BN(Q_j + \lambda \times M_j), \quad (5)$$

where  $j$  represents the class index of each occurring sound event that we can get from the clip-level labels of each audio clip, and  $\lambda$  denotes the update rate. Then, the enhancement of frame-level features can be formulated as:

$$X^e = X + X \odot Conv(Q_j), \quad (6)$$

**Table 2.** Performance comparison of CI-WSSSED and baseline models on the weakly labelled UrbanSED test set.

Method	AT-F1	Seg-F1	Event-F1
Base-CNN [17]	-	0.560	-
HTSAT [14]	0.771	0.644	0.210
CDur [13]	0.771	0.647	0.217
HTSAT-CI	<b>0.776</b>	<b>0.646</b>	<b>0.216</b>
CDur-CI	<b>0.774</b>	<b>0.648</b>	<b>0.220</b>

where  $\odot$  denotes matrix dot product and  $Conv$  represents the  $1 \times 1$  convolution. In this way, we can not only mitigate the impact of “entangled context” including entangled co-occurring sound events and background sounds, but also spotlight the active regions of the frame-level features, thus reducing classification errors and boosting localization performance.

During the training phase, our proposed network learns to minimize the cross-entropy losses for both classification branches. Specifically, we adopt two classifiers with shared weights for the two branches. The first classifier is used to produce initial prediction scores  $S = \{s_1, s_2, \dots, s_k\}$ , and the second classifier is accountable for generating more accurate prediction scores  $S^e = \{s_1^e, s_2^e, \dots, s_k^e\}$  using the enhanced frame-level features. Then, the cross-entropy losses for both classification branches are optimized to train the two classifiers together in an end-to-end pipeline. The overall loss function  $L$  is formulated below:

$$L = \left(-\sum_{i=1}^k s_i^* \log(s_i)\right) + \left(-\sum_{i=1}^k s_i^* \log(s_i^e)\right), \quad (7)$$

where  $s^*$  is the ground-truth label of an audio clip. While in the inference stage, we use the enhanced features to produce the final frame-level prediction results.

## 4. EXPERIMENTS AND RESULTS

### 4.1. Dataset

We evaluate our method on the two public sound event detection datasets: DCASE2017 task4 [3] and UrbanSED [17] datasets. The DCASE2017 task4 – Large-scale weakly supervised sound event detection for smart cars dataset is comprised of a training subset with 51172 audio clips, a validation subset with 488 audio clips, and an evaluation set with 1103 audio clips, including 17 sound events. The UrbanSED dataset has 10 event labels within urban environments, divided into 6000 training, 2000 validation, and 2000 evaluation audio clips.

### 4.2. Baseline Models and Training Details

To evaluate the effectiveness and generalization of our CI-WSSSED, we apply our method to multiple baseline systems,

including CDur [13], CNN-biGRU [1], CNN-Transformer [1] and HTSAT [14]. CDur consists of a 5-layer CNN and a bidirectional Gated Recurrent Unit (biGRU), while the CNN-biGRU system is modeled by a 9-layer CNN with a biGRU. The CNN-Transformer consists of a 9-layer CNN with one transformer block. HTSAT uses the Swin Transformer [18] backbone with ImageNet-pretraining, where we use 3 network groups with 2, 2, 6 swin-transformer blocks for the DCASE2017 task4 dataset, and for the UrbanSED dataset, we only adopt two stages with 2, 2 swin-transformer blocks. The update rate is set as  $\lambda = 0.01$ . We use audio tagging mAP (AT-mAP), audio tagging F1 score (AT-F1), sound event detection mAP (SED-mAP), Segment-F1 score (Seg-F1) and Event-F1 score to evaluate our method.

### 4.3. Results on DCASE2017 Task4

We first report our experiment results on both DCASE2017 task4 validation set and evaluation set in Table 1. We use baseline-CI to represent baseline models using our causal intervention method. It can be seen that our method achieves significant performance boosts on all baseline models, especially on the mAP, AT-F1 and Event-F1 metrics, which demonstrates the effectiveness of our CI-WSSSED in reducing classification errors and localizing entire sound events.

### 4.4. Results on UrbanSED

As shown in Table 2, we compare previous methods on the UrbanSED dataset with our CI-WSSSED. It is clear that our CI-WSSSED also achieves consistent improvements compared to the previous corresponding models. Notably, after applying the causal intervention method, the performance gain of UrbanSED is not as significant as that of the DCASE2017 task4 dataset. We infer the reason is that there are many sound event classes within the DCASE2017 task4 dataset that often co-occur [13], such as the sound events of “train” and “train horn”, as well as “car”, “car alarm”, and “car passing by”, so the DCASE2017 task4 dataset suffers more from the “entangled context” and thus benefits more from our CI-WSSSED.

## 5. CONCLUSIONS

In this paper, we target the “entangled context” problem in the WSSSED task from both aspects of entangled co-occurring sound events and background sounds. Through analyzing the causal relationship between frame-level features, contexts, and clip-level labels with the help of the SCM, we pinpoint the context as a co-occurrence confounder and then propose an end-to-end CI-WSSSED method to deal with the effect. Experiments show that the “entangled context” is a practical issue within the WSSSED task and our CI-WSSSED pipeline can effectively boost the performance of WSSSED on multiple datasets and generalize to various baseline models.

## 6. REFERENCES

- [1] Qiuqiang Kong, Yong Xu, Wenwu Wang, and Mark D Plumbley, "Sound event detection of weakly labelled data with cnn-transformer and automatic threshold optimization," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 2450–2460, 2020.
- [2] Yifei Xin, Dongchao Yang, and Yuexian Zou, "Audio pyramid transformer with domain adaption for weakly supervised sound event detection and audio classification," *Proc. Interspeech 2022*, pp. 1546–1550, 2022.
- [3] Annamaria Mesaros, Toni Heittola, Aleksandr Diment, Benjamin Elizalde, Ankit Shah, Emmanuel Vincent, Bhiksha Raj, and Tuomas Virtanen, "Dcase 2017 challenge setup: Tasks, datasets and baseline system," in *DCASE 2017-Workshop on Detection and Classification of Acoustic Scenes and Events*, 2017.
- [4] Keisuke Imoto and Seisuke Kyochi, "Sound event detection using graph laplacian regularization based on event co-occurrence," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 1–5.
- [5] Liwei Lin, Xiangdong Wang, Hong Liu, and Yueliang Qian, "Specialized decision surface and disentangled feature for weakly-supervised polyphonic sound event detection," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 1466–1478, 2020.
- [6] Tomer Galanti, Ofir Nabati, and Lior Wolf, "A critical view of the structural causal model," *arXiv preprint arXiv:2002.10007*, 2020.
- [7] Wei Qin, Hanwang Zhang, Richang Hong, Ee-Peng Lim, and Qianru Sun, "Causal interventional training for image recognition," *IEEE Transactions on Multimedia*, 2021.
- [8] Madelyn Glymour, Judea Pearl, and Nicholas P Jewell, *Causal inference in statistics: A primer*, John Wiley & Sons, 2016.
- [9] Dong Zhang, Hanwang Zhang, Jinhui Tang, Xian-Sheng Hua, and Qianru Sun, "Causal intervention for weakly-supervised semantic segmentation," *Advances in Neural Information Processing Systems*, vol. 33, pp. 655–666, 2020.
- [10] Thomas C Chalmers, Harry Smith Jr, Bradley Blackburn, Bernard Silverman, Biruta Schroeder, Dinah Reitman, and Alexander Ambroz, "A method for assessing the quality of a randomized control trial," *Controlled clinical trials*, vol. 2, no. 1, pp. 31–49, 1981.
- [11] Riddhiman Adib, Paul Griffin, Sheikh Iqbal Ahamed, and Mohammad Adibuzzaman, "A causally formulated hazard ratio estimation through backdoor adjustment on structural causal model," in *Machine Learning for Healthcare Conference*. PMLR, 2020, pp. 376–396.
- [12] Donmoon Lee, Subin Lee, Yoonchang Han, and Kyogu Lee, "Ensemble of convolutional neural networks for weakly-supervised sound event detection using multiple scale input," *Detection and classification of acoustic scenes and events (DCASE)*, 2017.
- [13] Heinrich Dinkel, Mengyue Wu, and Kai Yu, "Towards duration robust weakly supervised sound event detection," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 887–900, 2021.
- [14] Ke Chen, Xingjian Du, Bilei Zhu, Zejun Ma, Taylor Berg-Kirkpatrick, and Shlomo Dubnov, "Hts-at: A hierarchical token-semantic audio transformer for sound classification and detection," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 646–650.
- [15] Feifei Shao, Yawei Luo, Li Zhang, Lu Ye, Siliang Tang, Yi Yang, and Jun Xiao, "Improving weakly supervised object localization via causal intervention," in *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 3321–3329.
- [16] Guiguang Ding, Minghai Chen, Sicheng Zhao, Hui Chen, Jungong Han, and Qiang Liu, "Neural image caption generation with weighted training and reference," *Cognitive Computation*, vol. 11, no. 6, pp. 763–777, 2019.
- [17] Justin Salamon, Duncan MacConnell, Mark Cartwright, Peter Li, and Juan Pablo 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.
- [18] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10012–10022.