# AN EFFICIENT TEMPORARY DEEPFAKE LOCATION APPROACH BASED EMBEDDINGS FOR PARTIALLY SPOOFED AUDIO DETECTION

*Yuankun Xie, Haonan Cheng, Yutian Wang, Long Ye*\*

State Key Laboratory of Media Convergence and Communication,
Communication University of China, Beijing 100024, China

## ABSTRACT

Partially spoofed audio detection is a challenging task, lying in the need to accurately locate the authenticity of audio at the frame level. To address this issue, we propose a fine-grained partially spoofed audio detection method, namely Temporal Deepfake Location (TDL), which can effectively capture information of both features and locations. Specifically, our approach involves two novel parts: embedding similarity module and temporal convolution operation. To enhance the identification between the real and fake features, the embedding similarity module is designed to generate an embedding space that can separate the real frames from fake frames. To effectively concentrate on the position information, temporal convolution operation is proposed to calculate the frame-specific similarities among neighboring frames, and dynamically select informative neighbors to convolution. Extensive experiments show that our method outperform baseline models in ASVspoof2019 Partial Spoof dataset and demonstrate superior performance even in the cross-dataset scenario.

***Index Terms***— partially spoofed audio detection, temporal deepfake location, embedding learning.

## 1. INTRODUCTION

AI generated content (AIGC) technology has witnessed swift progress in recent years, particularly in speech-related applications like text-to-speech (TTS) [1, 2, 3] and voice conversion (VC) [4, 5, 6]. Although these technologies have brought about convenience, they have also posed significant security threats. Thus, various initiatives and challenges, such as ASVspoof [7, 8], have been established to foster research on countermeasure solutions that safeguard speech applications and human listeners against spoofing attacks. Nevertheless, a significant scenario has been overlooked in most datasets and challenges where a bonafide speech utterance is contaminated by synthesized speech segments, leading to partial spoofing (PS). Attackers can use PS to alter sentence semantics, and such modifications can be easily accomplished at low cost. For instance, attackers can easily modify single word such as time, place, and characters in sentence to dramatically change the semantics. Furthermore, If attackers have knowledge of phonology, they can manipulate vowels and even consonants such as "pan,""pin,""pen," which are smaller than the word level. Therefore, defending against such fine-grained PS scenarios poses significant challenges for defenders.

In recent years, there are several studies about PS scenarios for Audio Deepfake Detection (ADD). Yi et al. [9] create a dataset that focuses on changing a few words in an utterance for half-truth audio

detection. At the same time, Zhang et al. [10] construct a speech database called 'PartialSpoof' designed for PS scenarios. The above two datasets are the beginning of the research for PS scenario in ADD task. Afterward, Zhang et al. [11] propose the SELCNN network to enhance the ability of the accuracy of the utterance. Lv et al. [12] use Wav2Vec2 (W2V2) [13] as front-end, ECAPA-TDNN [14] as back-end achieving the first rank in ADD 2022 Track 2[15]. Although the above research shows effectiveness at the utterance level detection in PS, they do not pinpoint specific segments with precision. Thus, Zhang et al. [16] extended the previous utterance-level PS dataset labels to frame-level and proposed corresponding W2V2-based countermeasures to enhance frame-level detection capability.

The aforementioned methods solely utilize existing ADD models such as LCNN, currently lacking specific approaches tailored to the PS scenario, particularly in terms of precise frame-level localization. To address this challenge, we propose a novel Temporal Deepfake Location (TDL) method. For front-end, we take advantage of W2V2 [17]. By training on a vast corpus of genuine speech from diverse source domains, W2V2 can effectively discriminate the real and fake in complex acoustic scenarios. For back-end, our primary focus is on fine-grained locating the genuine and spoofed speech segment. To clearly distinguish the real and fake in feature level, we first design the embedding similarity module to separate the real and fake frames in embedding space and get a high-quality embedding similarity vector. Then, we propose temporal convolution operation to locate the region from the embedding vector. The local similarity for each temporal position is calculated from the embedding. By this means, we can obtain a frame-specific weight to guide the convolution making a temporal sensitive calculation. Our main contributions can be summarized as follows:

- We propose TDL method, an efficient and effective ADD method for PS scenarios which combines a embedding similarity module and temporal convolution operation to effectively capture both feature and positional information.

- The proposed method outperforms baseline models in ASV spoof 2019PS dataset and demonstrate superior performance even in cross-dataset experiments.

## 2. PROPOSED METHOD

### 2.1. Problem statement and overview

In PS scenarios, the fake audio segment is inserted within the genuine speech. Our target is to detect the real and fake segments at frame level. Given the large-scale self-supervised audio feature $f = (f_1, f_2, ... f_T) \in R^{D \times T}$, where $D$ and $T$ denote the dimension of audio feature and the number of frames respectively. The whole
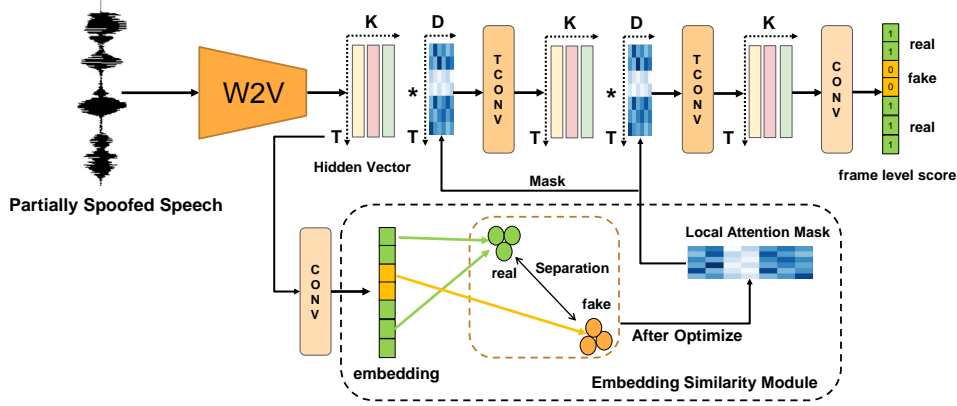
---

**Fig. 1**. The entire structure of our proposed Temporal Deepfake Location (TDL) method.

task is defined as input feature $f$ and output the frame level label $y = (y_1, y_2, ...y_T) \in \{0, 1\}^T$, where 1 represents the real frames and 0 represents the fake frames.

The framework of our proposed TDL is depicted in Figure 1. First, we utilize Wav2Vec-XLS-R to extract the frame level feature from the raw audio. Then, for enhanced identification of genuine and fake distinctions at the embedding level, we devise an embedding similarity module to segregate authentic and synthetic frames within the embedding space. Next, to capture the position information, we adopt temporal convolution operation by focusing on frame-specific similarities among neighboring frames. Finally, we employ 1D convolutional layers and fully connected layers for downsampling to the frame level label to compute the Binary Cross-Entropy (BCE).

### 2.2. W2V2 front-end

W2V2 based front-end is trained by solving a contrastive task over a masked feature encoder. Firstly, speech signals in various lengths are passed through a feature extractor consisting of seven convolutional neural network (CNN) layers. Subsequently, context representations are obtained using a Transformer network [18] comprising of 24 layers, 16 attention heads, and an embedding size of 1024. In practice, we utilize the Hugging Face version of wav2vec2-XLS-R-300M[1] and freeze the weights of the front-end. The front-end model is pre-trained with 436k hours of unannotated genuine speech data in 128 languages. Consequently, the last hidden states from the transformer can effectively represent the contextualized information of genuine speech which is different from the partially fake speech.

### 2.3. Embedding similarity module

To better capture feature-level information, we first distinguish the real and fake frames in the embedding space. Specifically, the W2V2 features are fed into a CONV module, consisting of two sequential 1D-CNNs, which downsamples the embedding dimension from 1024 to 32. The embedding vector is $L2$-normalized. Then we get a embedding vector $e = (e_1, e_2, ...e_T) \in R^{D \times T}$. In the embedding similarity module, we utilize cosine similarity to measure the similarity of two embedding vector $e_u$ and $e_v$ as follows:

$$\mathcal{S}(\mathbf{e}_u, \mathbf{e}_v) = \frac{\mathbf{e}_u^T \cdot \mathbf{e}_v}{\|\mathbf{e}_u\|_2 \cdot \|\mathbf{e}_v\|_2}. \qquad (1)$$

To increase the distance between genuine and fake frames in the embedding space and improve generalizability, we computed the cosine similarities between genuine frames, between fake frames, and between genuine and fake frames. Specifically, we ensured that genuine frames from different positions exhibited similarity, fake frames from different positions exhibited similarity, while genuine and fake frames are dissimilar to each other.

Thus, $\mathcal{L}_{ESM}^{Real}$ and $\mathcal{L}_{ESM}^{Fake}$ are proposed to make the real frames and fake frames in different positions similar:

$$\mathcal{L}_{\text{ESM}}^{\text{Real}} = \max_{\forall e_x, e_y, x \neq y} \lfloor \tau_{\text{same}} - \mathcal{S}(\mathbf{e}_x, \mathbf{e}_y) \rfloor_+, \qquad (2)$$

$$\mathcal{L}_{\text{ESM}}^{\text{Fake}} = \max_{\forall e_m, e_n, m \neq n} \lfloor \tau_{\text{same}} - \mathcal{S}(\mathbf{e}_m, \mathbf{e}_n) \rfloor_+, \qquad (3)$$

where $e_x$ and $e_y$ refer to distinct positions of real frames, while $e_m$ and $e_n$ refer to those of fake frames. $\tau_{\text{same}}$ is the similarity threshold between frames from the same category, $\lfloor ... \rfloor_+$ represents clipping below at zero. It is noteworthy that although we know the positions of frame-level authenticity labels, the temporal dimension of W2V2-XLS-R features does not inherently align with these frame-level labels. To tackle this issue, we ascertain the temporal authenticity in the time dimension of the embedding vector by calculating the ratio between the temporal dimensions of the label and the embedding vector.

$\mathcal{L}_{ESM}^{Diff}$ is proposed to separate the real and fake frames, which can be formulated as:

$$\mathcal{L}_{\text{ESM}}^{\text{Diff}} = \max_{\forall e_r, e_f} \lfloor \mathcal{S}(\mathbf{e}_r, \mathbf{e}_f) - \tau_{\text{Diff}} \rfloor_+, \qquad (4)$$

where $e_r$ and $e_f$ refer to the embedding vector of real frames and fake frames. $\tau_{\text{diff}}$ is the similarity threshold to constrain the distance between real and fake frames. Finally, the embedding similarity module is optimized by $\mathcal{L}_{ESM}$, which takes into account the three aforementioned losses in a joint manner. The $\mathcal{L}_{ESM}$ is calculated as follows:

$$\mathcal{L}_{ESM} = \mathcal{L}_{ESM}^{Real} + \mathcal{L}_{ESM}^{Fake} + \mathcal{L}_{ESM}^{Diff}. \qquad (5)$$

### 2.4. Temporal convolution operation

To effectively capture the positional information, we use the embedding vector as an local attention mask to perform temporal convolution operations. Consider a audio feature $\mathbf{X} \in R^{D_{in} \times T}$, where $D_{in}$ and $T$ represent the dimension of the vector and number of frames

**Table 1**. Architecture of TDL network.

| module | kernel/stride | output shape |
|---|---|---|
| W2V2 | - | (batch,1024,1050) |
| CONV | 3/1<br>3/1 | (batch,512,1050)<br>(batch,32,1050) |
| TCONV | 3/1 | (batch,1024,1050) |
| TCONV | 3/1 | (batch,1024,1050) |
| CONV | 1/1 | (batch,2,1050) |
| Flatten/FC | - | (batch,132) |

respectively. The temporal convolution layer learns a dynamic convolution kernel $\Bbbk \in R^{k \times D_{in} \times D_{out}}$, where $k$ is the size of temporal kernel, $D_{out}$ is the dimension of output feature. We only utilize the dynamic kernel $\Bbbk^m \in R^{k \times D_{in}}$ to compute $m^{th}$ channel of the output for convenient. Thus, the temporal convolution operation for the $t^{th}$ feature can be expressed as:

$$f_t^m = \sum_{i=0}^{k-1} \Bbbk^m[i,:] \cdot \overline{\mathbf{X}}\left[:, t - \frac{k}{2} + i\right], \tag{6}$$

where $f_t^m$ is the value in the $m^{th}$ channel of output feature vector, $[\cdots]$ means a slice of a matrix, $(\cdot)$ denotes the inner product. $\overline{\mathbf{X}}$ is the modulated feature processed by neighbor similarity calculation:

$$\overline{\mathbf{X}}\left[:, t - \frac{k}{2} + i\right] = \mathbf{X}\left[:, t - \frac{k}{2} + i\right] \times \mathbf{a}[i,t], \tag{7}$$
$$i \in [0, \ldots, k-1],$$

where matrix $\mathbf{a} \in R^{k \times T}$ is a similarity matrix that calculate the local similarity for each temporal position, $\mathbf{a}[i,t]$ indicates the similarity between the $t^{th}$ feature vector and its $k$ neighbors.

In practice, we determine the dynamic kernel weight based on the embedding vector generated by ESM module. We apply temporal convolution operation to the W2V2 features on two sequence 1D-CNNs, where both input channel and output channel remain unchanged to maintain consistency in temporal dimension.

### 2.5. Total loss

Following two consecutive temporal convolution operation layers, to capture additional temporal information and align with the label dimensions, we subsequently employ 1D-CNN, fully connected (FC) layers, and sigmoid activation functions to calculate the BCE loss. The architecture details of TDL is shown in Table 1. The total loss is defined as follow:

$$\mathcal{L}_{all} = \mathcal{L}_{BCE} + \lambda\mathcal{L}_{ESM}, \tag{8}$$

where $\lambda$ is set to 0.1 to balance the value of two losses.

## 3. EXPERIMENTS

### 3.1. Database

Our experiments for PS scenario include two public datasets: ASVspoof2019PS (19PS) [10] and LAV-DF [19]. 19PS is constructed based on the ASVspoof2019 LA database [20]. All experiments on the 19PS dataset are conducted using 160ms resolution labels. The training, validation, and testing sets are distributed according to the original dataset allocation, consisting of 25,380, 24,844, and 71,237 utterances respectively.

**Table 2**. Percentages(%) of fake class in each dataset.

| dataset | subset | frame-level | utterance-level |
|---|---|---|---|
| 19PS | train | 53.00 | 89.83 |
| 19PS | dev | 52.31 | 89.74 |
| 19PS | test | 48.03 | 89.68 |
| LAV-DF | test | 10.01 | 48.82 |

To evaluate the model's generalizability, we conduct additional testing of the 19PS-trained model using the LAV-DF test set. LAV-DF represents a multi-modal temporal forgery dataset, containing a total of 26,100 videos in test set. We extract the audio track of each video and create 160ms frame level genuine and fake labels.

We calculated the percentage of samples belonging to fake class at both the frame and sentence levels, as shown in the Table 2. We can observe that the frame-level labels in 19PS are balanced, facilitating model training. However, the LAV-DF dataset exhibits a lower proportion of spoof segments, making it unbalanced and presenting greater challenges for detection.

### 3.2. Implementation details

In order to address the issue of variable-length audio inputs, we employ the technique of zero-padding to the maximum length of training set. For the frame of genuine speech, we set the label to one, while for spoofing frame, the label is set to zero. In the case of 19PS, the maximum duration of speech in the training set is 21.03 seconds with a W2V2 feature dimension of (1050,1024) and the number of frames at a resolution of 160 ms is 132. For LFCC, we extracted 60-dimensional LFCC with a combination of static, delta and delta-delta coefficients.

For training strategy, the Adam optimizer is adopted with $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\varepsilon = 10^{-9}$ and weight decay is $10^{-4}$. We train all of the models for 100 epochs. The learning rate is initialized as $10^{-5}$ and halved every 5 epochs. It is worth mention that no data augmentation method is used for experiment.

### 3.3. Evaluation metrics

In our experiment, we employ four evaluation metrics to assess model performance: Equal error rate (EER), precision, recall, and $F_1$ score. All metrics are computed based on frame-level authenticity labels of the partially spoofed audio. Precision, recall, and $F_1$ score are defined as follow:

$$Precison = \frac{TP}{TP + FP}, \tag{9}$$

$$Recall = \frac{TP}{TP + FN}, \tag{10}$$

$$F_1 score = \frac{2 \cdot Precison \cdot Recall}{Precison + Recall}, \tag{11}$$

where $TP$, $TN$, $FP$, $FN$ represent the numbers of true positive, true negative, false positive, and false negative samples, respectively. In practice, we employed point-based binary-classification precision, recall, and $F_1$ score from Sklearn. Before any evaluation, zero-padding is eliminated based on the actual length of the features.

**Table 3**. EER results (%) on ASVspoof2019 PS dataset.

| Model | Feature | EER |
|---|---|---|
| LCNN-BLSTM [10] | LFCC | 16.21 |
| SELCNN-BLSTM [11] | LFCC | 15.93 |
| LCNN-BLSTM [10] | W2V2-XLS-R | 9.87 |
| 5gMLP [16] | W2V2-Large | 9.24 |
| TDL (w/o ESM) | W2V2-XLS-R | 8.79 |
| TDL | W2V2-XLS-R | **7.04** |

**Table 4**. The four evaluation metrics results (%) for training on 19PS and testing on the LAV-DF test set.

| Model | Feature | EER↓ | Precision↑ | Recall↑ | $F_1$ score↑ |
|---|---|---|---|---|---|
| LCNN-BLSTM | LFCC | 17.89 | 95.93 | 73.73 | 83.38 |
| LCNN-BLSTM | W2V2-XLS-R | 15.35 | **99.05** | 62.32 | 76.50 |
| TDL | W2V2-XLS-R | **11.23** | 98.73 | **75.42** | **85.51** |

## 4. RESULTS AND DISCUSSIONS

### 4.1. Results

**Results on 19PS.** We compare the performance of several baseline models in terms of EER metric, as presented in Table 3. All models are trained on the 19PS training dataset. TDL (w/o ESM) represents our model without ESM module. As shown in Table 3, our model achieve the lowest EER 7.04% in partially spoofed audio detection task.

Based on the experimental results, We first observe that the impact of feature is greater than backbone. For instance, as seen in first and third row in Table 3, where the backbone is LCNN-BLSTM, the utilization of W2V2 features resulted in a 6.84% EER decrease compared to LFCC. Conversely, when feature remain consistent, as demonstrated in first and second row of the in Table 3, both employing the shared LFCC attribute, SELCNN-LSTM exhibited a marginal EER reduction of 0.28% in comparison to LCNN-LSTM. Furthermore, we find that the architecture design of the TDL network aligns well with partial spoofed detection. Specifically, when the features utilized W2V2-XLS-R, the TDL (without ESM module) still exhibits a 1.08% reduction in EER compared to the LCNN-BLSTM.

**Results on LAV-DF.** To validate the generalizability of our proposed model, we train on 19PS and evaluate on the test set of LAV-DF for 4 evaluation metrics. The results of the testing are presented in the Table 4. Although LAV-DF is an unbalanced dataset, our proposed model achieve the best performance of 11.23% EER compared to baseline models.

### 4.2. Disscussion

**Label Setting.** As we mentioned in Section 3.2, we set real frames, fake frames for 1 and 0. To the best of our knowledge, there has been no prior research discussing which label configuration will be beneficial to the final prediction. Therefore, we experiment with three different label settings on our proposed TDL model as shown in Table 5.

"Boundary 1" indicates that we set the boundary frames between genuine and fake segments as 1, while other positions are set as 0. In practice, due to the sparsity of boundary frames, we set 4 boundary frames at the transition between genuine and fake segments. Addi-

**Table 5**. EER results (%) of different label configuration on 19PS.

| Label Setting | EER↓ | Precision↑ | Recall↑ | $F_1$ score↑ |
|---|---|---|---|---|
| Boundary 1 | 10.89 | 79.72 | 82.01 | 80.85 |
| real 0 fake 1 | 9.52 | 81.87 | 84.52 | 83.17 |
| real 1 fake 0 | **7.04** | **88.69** | **95.01** | **91.54** |

**Table 6**. Parameters (in thousands) comparison.

| Model | Parameters |
|---|---|
| TDL | 8,718 |
| LCNN-BLSTM | 21,511 |

tionally, we employ a weighted BCE loss, assigning a weight value of 100 to the boundary values, as a replacement for standard BCE. Experimental results demonstrate that this method is less effective compared to directly predicting the authenticity of individual frames. Additionally, since predicting boundaries often requires further verification of the genuineness of the segments on both sides, we did not adopt the boundary setting.

For the frame-level direct prediction of authenticity, we conducted experiments by setting real frames as 0 and fake frames as 1, and alternatively by setting real frames as 1 and fake frames as 0, as shown in the "real 0 fake 1" and "real 1 fake 0" of the Table 5 respectively. Experiments results show that "real 1 fake 0" outperform "real 0 fake 1" in four evaluation metrics, especially in recall metric, which indicates that TDL can accurately identify genuine speech. When setting real frames as "1" and fake frames along with padding frames as "0", we can better concentrate on the real segment. This is similar to previous works [21, 22] which also focus on the real speech distribution in fully-spoofed ADD task. Through our experiments, we have demonstrated that it is also significant in partially-spoofed ADD task. This is also why W2V2 features are effective in the field of ADD which only extracted by rich real source domains.

**Complexity Comparision.** Apart from evaluating the performance, we measured the complexity of the models. For frame-level detection task, particularly for fine-grained prediction, the large final output dimension can result in excessive parameterization and low efficiency. Unlike LCNN, which convolves overall values, our proposed TDL model uses temporal convolution operation to selectively focus only on high-weight regions. It can be observed that the parameter count of TDL is only 40.53% of that of LCNN-BLSTM, which is shown in Table 6.

## 5. CONCLUSION

In this paper, we propose an efficient temporary deepfake location approach based embeddings for partially spoofed audio detection. TDL can achieve outstanding performance benefits from two designed core modules: embedding similarity module and temporal convolution operation, which can effectively capture both feature and positional information. The experimental results demonstrate that TDL achieves the best performance in the 19PS dataset and also perform well in cross-dataset scenarios.

## 6. REFERENCES

[1] Rongjie Huang, Zhou Zhao, Huadai Liu, Jinglin Liu, Chenye Cui, and Yi Ren, "Prodiff: Progressive fast diffusion model

for high-quality text-to-speech," in *Proceedings of ACM MM*, 2022, pp. 2595–2605.

[2] X. Tan, J. Chen, H. Liu, J. Cong, C. Zhang, Y. Liu, X. Wang, Y. Leng, Y. Yi, L. He, et al., "Naturalspeech: End-to-end text to speech synthesis with human-level quality," *arXiv preprint arXiv:2205.04421*, 2022.

[3] Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang, Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming Wang, Jinyu Li, et al., "Neural codec language models are zero-shot text to speech synthesizers," *arXiv preprint arXiv:2301.02111*, 2023.

[4] Chak Ho Chan, Kaizhi Qian, Yang Zhang, and Mark Hasegawa-Johnson, "Speechsplit2. 0: Unsupervised speech disentanglement for voice conversion without tuning autoencoder bottlenecks," in *Proceedings of ICASSP*, 2022, pp. 6332–6336.

[5] Y. Chen, D. Wu, T. Wu, and H. Lee, "Again-vc: A one-shot voice conversion using activation guidance and adaptive instance normalization," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 5954–5958.

[6] Huaizhen Tang, Xulong Zhang, Jianzong Wang, Ning Cheng, and Jing Xiao, "Avqvc: One-shot voice conversion by vector quantization with applying contrastive learning," in *Proceedings of ICASSP*. IEEE, 2022, pp. 4613–4617.

[7] Andreas Nautsch, Xin Wang, Nicholas Evans, Tomi H Kinnunen, Ville Vestman, Massimiliano Todisco, Héctor Delgado, Md Sahidullah, Junichi Yamagishi, and Kong Aik Lee, "Asvspoof 2019: spoofing countermeasures for the detection of synthesized, converted and replayed speech," *IEEE Transactions on Biometrics, Behavior, and Identity Science*, vol. 3, no. 2, pp. 252–265, 2021.

[8] H. Delgado, N. Evans, T. Kinnunen, K. Lee, X. Liu, A. Nautsch, J. Patino, M. Sahidullah, M. Todisco, X. Wang, et al., "Asvspoof 2021: Automatic speaker verification spoofing and countermeasures challenge evaluation plan," *arXiv preprint arXiv:2109.00535*, 2021.

[9] Jiangyan Yi, Ye Bai, Jianhua Tao, Haoxin Ma, Zhengkun Tian, Chenglong Wang, Tao Wang, and Ruibo Fu, "Half-truth: A partially fake audio detection dataset," in *Proceedings of Interspeech*, 2021, pp. 1654–1658.

[10] Lin Zhang, Xin Wang, Erica Cooper, Junichi Yamagishi, Jose Patino, and Nicholas Evans, "An initial investigation for detecting partially spoofed audio," in *Proceedings of Interspeech*, 2021, pp. 4264–4268.

[11] Lin Zhang, Xin Wang, Erica Cooper, and Junichi Yamagishi, "Multi-task learning in utterance-level and segmental-level spoof detection," *arXiv preprint arXiv:2107.14132*, 2021.

[12] Zhiqiang Lv, Shanshan Zhang, Kai Tang, and Pengfei Hu, "Fake audio detection based on unsupervised pretraining models," in *Proceedings of ICASSP*, 2022, pp. 9231–9235.

[13] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in neural information processing systems*, vol. 33, pp. 12449–12460, 2020.

[14] Brecht Desplanques, Jenthe Thienpondt, and Kris Demuynck, "Ecapa-tdnn: Emphasized channel attention, propagation and aggregation in tdnn based speaker verification," *arXiv preprint arXiv:2005.07143*, 2020.

[15] Jiangyan Yi, Ruibo Fu, Jianhua Tao, Shuai Nie, Haoxin Ma, Chenglong Wang, Tao Wang, Zhengkun Tian, Ye Bai, Cunhang Fan, et al., "Add 2022: the first audio deep synthesis detection challenge," in *Proceedings of ICASSP*. IEEE, 2022, pp. 9216–9220.

[16] Lin Zhang, Xin Wang, Erica Cooper, Nicholas Evans, and Junichi Yamagishi, "The partialspoof database and countermeasures for the detection of short fake speech segments embedded in an utterance," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2022.

[17] Arun Babu, Changhan Wang, Andros Tjandra, Kushal Lakhotia, Qiantong Xu, Naman Goyal, Kritika Singh, Patrick von Platen, Yatharth Saraf, Juan Pino, et al., "Xls-r: Self-supervised cross-lingual speech representation learning at scale," *arXiv preprint arXiv:2111.09296*, 2021.

[18] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[19] Zhixi Cai, Kalin Stefanov, Abhinav Dhall, and Munawar Hayat, "Do you really mean that? content driven audio-visual deepfake dataset and multimodal method for temporal forgery localization," in *2022 International Conference on Digital Image Computing: Techniques and Applications (DICTA)*. IEEE, 2022, pp. 1–10.

[20] Xin Wang, Junichi Yamagishi, Massimiliano Todisco, Héctor Delgado, Andreas Nautsch, Nicholas Evans, Md Sahidullah, Ville Vestman, Tomi Kinnunen, Kong Aik Lee, et al., "Asvspoof 2019: A large-scale public database of synthesized, converted and replayed speech," *Computer Speech & Language*, vol. 64, pp. 101114, 2020.

[21] Y. Zhang, F. Jiang, and Z. Duan, "One-class learning towards synthetic voice spoofing detection," *IEEE Signal Processing Letters*, vol. 28, pp. 937–941, 2021.

[22] Yuankun Xie, Haonan Cheng, Yutian Wang, and Long Ye, "Learning A Self-Supervised Domain-Invariant Feature Representation for Generalized Audio Deepfake Detection," in *Proc. INTERSPEECH 2023*, 2023, pp. 2808–2812.