# Detect All-Type Deepfake Audio: Wavelet Prompt Tuning for Enhanced Auditory Perception

Yuankun Xie
xieyuankun@cuc.edu.cn
State Key Laboratory of Media
Convergence and Communication,
Communication University of China
Beijing, China

Ruibo Fu
Institute of Automation, Chinese
Academy of Sciences
Beijing, China

Zhiyong Wang
School of Artificial Intelligence,
University of Chinese Academy of
Sciences
Beijing, China

Xiaopeng Wang
School of Artificial Intelligence,
University of Chinese Academy of
Sciences
Beijing, China

Songjun Cao
Youtu Lab, Tencent
Beijing, China

Long Ma
Youtu Lab, Tencent
Beijing, China

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

Long Ye
yelong@cuc.edu.cn
State Key Laboratory of Media
Convergence and Communication,
Communication University of China
Beijing, China

## ABSTRACT

The rapid advancement of audio generation technologies has escalated the risks of malicious deepfake audio across speech, sound, singing voice, and music, threatening multimedia security and trust. While existing countermeasures (CMs) perform well in single-type audio deepfake detection (ADD), their performance declines in cross-type scenarios. This paper is dedicated to studying the all-type ADD task. We are the first to comprehensively establish an all-type ADD benchmark to evaluate current CMs, incorporating cross-type deepfake detection across speech, sound, singing voice, and music. Then, we introduce the prompt tuning self-supervised learning (PT-SSL) training paradigm, which optimizes SSL front-end by learning specialized prompt tokens for ADD, requiring 458× fewer trainable parameters than fine-tuning (FT). Considering the auditory perception of different audio types, we propose the wavelet prompt tuning (WPT)-SSL method to capture type-invariant auditory deepfake information from the frequency domain without requiring additional training parameters, thereby enhancing performance over FT in the all-type ADD task. To achieve an universally CM, we utilize all types of deepfake audio for co-training. Experimental results demonstrate that WPT-XLSR-AASIST achieved the best performance, with an average EER of 3.58% across all evaluation sets. The code is online available [1].

## KEYWORDS

Audio Deepfake Detection, Countermeasures, Prompt Tuning, Discrete Wavelet Transform
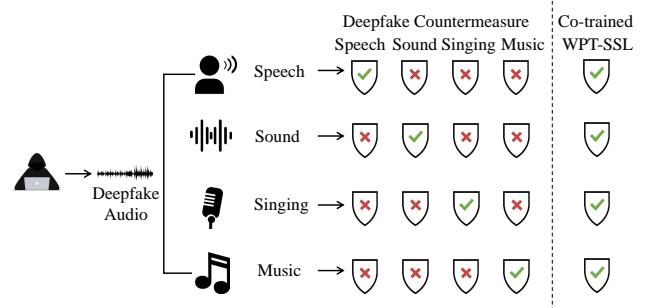
---

[1] https://github.com/xieyuankun/All-Type-ADD



**Figure 1: The challenge for current single-type trained CMs toward cross-type ADD task, highlighting the effectiveness of our proposed WPT-SSL CM.**

## 1 INTRODUCTION

With the development of audio language model (ALM) technology, it has become increasingly easy to synthesize any type of audio, including deepfake speech, sound, singing voice, and music. These deepfake audios pose a threat to society in various fields such as media, entertainment, cybersecurity, and political communication. Fortunately, research on audio deepfake detection (ADD) has been increasing annually. Among these, the earliest studies focused on deepfake speech detection. Researchers have developed a series of deepfake countermeasures (CMs) aimed at effectively detecting deepfake speech, based on the ASVspoof challenges [1–3]. Currently, some ADD research has gone beyond speech, such as the detection of deepfake singing voices [4–6], sounds [7, 8], and music [9].

Although each type of deepfake audio has its corresponding countermeasure (CM), in real-world scenarios, the type of audio is often uncertain and may encompass one or more categories, such as speech, sound, singing voice, or music. This leads to the challenge that the CM trained on a single type being unable to generalize and detect all types of audio, as shown in Figure 1. Therefore, it is crucial to develop an advanced CM that can generalize and effectively detect all-type of deepfake audio.

For CMs, the most effective approach currently is to use pre-trained self-supervised learning (SSL) features along with a classification backbone. A representative CM in speech deepfake detection is XLSR-AASIST [10], which fine-tunes (FT) the wav2vec2-xls-r (XLSR) [11] model on speech deepfake detection dataset, achieving excellent intra-domain (ID) and out-of-domain (OOD) generalization performance. However, when dealing with the all-type ADD task, several challenges are encountered. Firstly, from the data perspective, it is uncertain whether a CM trained on single audio type can generalize to detect other types of deepfake audio. Although some studies have investigated cross-type detection for two types [8, 12], there has been no exploration of cross-type detection for all audio types. Secondly, there has been no investigation into whether a domain-invariant feature exists that can ensure the invariance of authenticity discrimination across different audio types. This requires a detailed investigation of various SSL features as well as handcrafted features. Lastly, concerning the algorithm, although fine-tuning can yield promising results, it is highly dependent on specific hyper-parameters and requires a significant amount of training parameters [13].

To address the aforementioned challenges, in this paper, we aim to develop an all-type audio deepfake CM. We are the first to comprehensively establish an all-type ADD benchmark, which includes cross-type deepfake detection among speech, sound, singing voice, and music. For the feature of CMs, we investigate handcrafted features, raw waveforms, and various SSL-based features through both freezing and fine-tuning. For the classifier, we use AASIST [14], the most popular model in the field of ADD, as the back-end, and combine it with SSL front-end to form SSL-AASIST.

To efficiently optimize SSL front-end, inspired by Visual Prompt Tuning (VPT) [15], we proposed the Prompt Tuning (PT)-SSL training paradigm for ADD task. PT-SSL introduces learnable prompt tokens before the input of each transformer layer, while keeping the other parameters of the layers frozen, with the goal of learning specialized prompt tokens for the ADD task. Furthermore, considering the human cognition of different audio types, the primary differences in perceiving audio types lie in their frequency domain distributions [16–18]. However, current SSL models like wav2vec2, which are primarily designed for speech recognition, focus on temporal and specific speech frequency information, lacking the ability to capture full-frequency information. To enhance frequency domain adaptability and enable SSL-based CM to quickly adapt to all types of deepfake audio, we propose wavelet prompt learning (WPT)-SSL method. WPT-SSL applies a discrete wavelet transform (DWT) to a portion of the prompt tokens, obtaining tokens for different frequency bands, thereby enhancing the full-frequency perception capability of SSL-based CM. Surprisingly, we discovered that WPT-SSL can learn a type-invariant deepfake detection

prompt in a specific frequency band (HH) obtained through wavelet decomposition, thereby enabling all-type audio deepfake detection.

We summarize the contributions of this work as follow:

- We proposed all-type ADD task and established a comprehensive benchmark to measure the current CM's capability in detecting all-type deepfake audio.
- To efficiently train SSL front-end, we proposed the PT-SSL training paradigm, which significantly reduces the number of training parameters by only learning prompt tokens, achieving performance close to FT.
- Considering the human perception of different audio types, we proposed the WPT-SSL method, which can learn type-invariant frequency authenticity information. Without adding extra training parameters, WPT outperformed FT under all ADD test conditions.
- To achieve an universally CM, we utilize all types of deepfake audio for co-training. Experimental results demonstrate that WPT-SSL-AASIST achieved the best performance with an average EER of 3.58%.
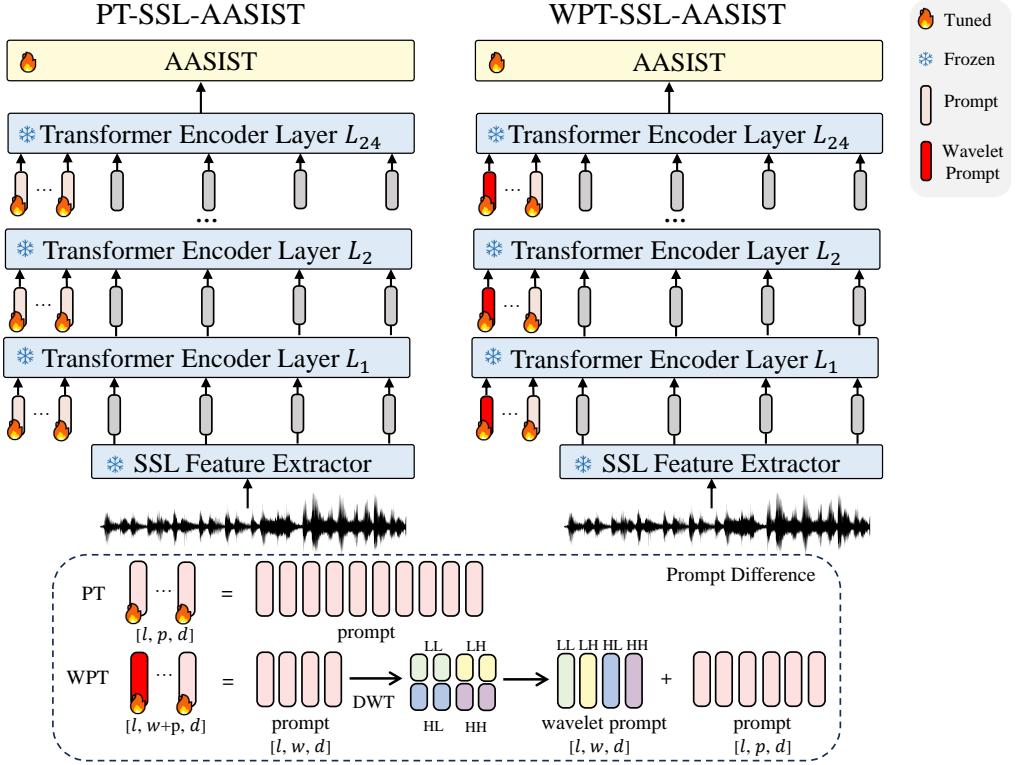
## 2 RELATED WORK

In this section, we primarily review the current popular ADD algorithms and the prompt tuning-related algorithms.

### 2.1 Audio Deepfake Detection

In introducing the current work in the field of ADD, we categorize it based on types: speech, sound, singing voice, and music.

**Speech**: Speech deepfake detection is the most extensively researched and thoroughly studied task among the four types. Surrounding the ASVspoof series, a number of outstanding studies have been developed. For CM without an SSL front-end, Jung et al. [14] proposed AASIST, which achieved an EER of 0.83% on the ASVspoof2019LA (19LA) test set for the first time without SSL features, by employing raw waveform input and a spectro-temporal graph attention method. For CM with SSL front-end, Tak et al. [10] proposed using XLSR as a front-end for fine-tuning, followed by the AASIST back-end, achieving milestone scores in the ADD field on the ASVspoof2021LA (21LA) and ASVspoof2021DF (21DF) datasets. Subsequent research has focused on exploring different categories of SSL [19, 20], the utilization of SSL layers [13, 21, 22], and robustness of CM [23, 24].

**Other types (sound, singing, music)**: Compared to the work on speech deepfake detection, research on other types is still primarily at the dataset stage, with fewer methodological explorations. For sound, Xie et al. [8] introduced the A3 subset of the Codecfake dataset, which includes deepfake sounds generated by Audiogen [25]. Furthermore, Xie et al. introduced the FakeSound dataset [7], developing novel countermeasures for sound deepfake detection. In the realm of singing voice, numerous advanced studies [26–28] have been proposed in relation to the SVDD challenge [4]. These studies indicate that, in addition to XLSR, the integration of other SSL models such as MERT [29] and WavLM [30] can further enhance CM performance. In the field of music, Luca et al. proposed FakeMusicCaps [9], which covers various ALM music synthesis methods. Concurrently, there have been limit studies on fake music detection [31, 32].

**Figure 2: Our proposed PT-SSL-AASIST (left) and WPT-SSL-AASIST (right). The differences between PT and WPT are illustrated below.**

**Cross types**: Regarding cross-type ADD task, some studies have explored the deepfake detection effects between two audio types. For instance, Gohari et al. [12] investigated cross-type ADD between speech and singing voice, while Li et al. [33] investigated the extension of speech deepfake detection to music. However, there is currently no research covering all-audio types.

## 2.2 Prompt Tuning

VPT [15] was the first to successfully introduce the PT [34] method into the visual domain. It consists of the VPT-SHALLOW version, which incorporates learnable prompts before the first Transformer layer of VIT, and the VPT-DEEP version, which prepends learnable parameters to the input of each Transformer encoder layer. In VTAB-1k [35], VPT-DEEP surpasses most fine-tuning methods and conserves a substantial number of parameters. The VPT approach and its enhancements have shown promising results in multiple downstream tasks [36–38].

Owing to the strong adaptability VPT to downstream tasks and minimal storage parameters, audio researchers have begun to apply the VPT paradigm to audio tasks [39, 40]. Recently, Oiso et al. [41] investigated the use of the PT paradigm for ADD tasks with limited data, and their study was the first to demonstrate the effectiveness of the PT paradigm in the ADD domain. However, they only utilized the SHALLOW PT paradigm with the first layer combined with prompts and tuned using target domain data. The potential of

PT paradigm in speech deepfake detection and all-type deepfake detection has yet to be fully explored.

## 3 PROMPT TUNING COUNTERMEASURE

In this section, we introduce our proposed PT-SSL-AASIST and WPT-SSL-AASIST paradigm, which rapidly adapt SSL features to the ADD task by learning prompt tokens.

## 3.1 PT-SSL-AASIST

For an input audio $X$, we first pad or chop it to a fixed length $L$, obtaining the audio input $X \in \mathbb{R}^L$. Then, the audio input is first passed through the frozen SSL front-end feature extractor. For SSL implementations such as XLSR, this feature extractor comprises a 7-layer CNNs. Subsequently, we obtain the input to the first encoder layer of the transformer, $E_0 \in \mathbb{R}^{t \times d}$, where $t$ represents the temporal length of the audio sequence and $d$ denotes the dimension of the transformer hidden states. For the prompt token, we employ Xavier uniform initialization for all layers, resulting in $\mathbf{P} = \left\{ \mathbf{P}_k \in \mathbb{R}^{p \times d} \mid k \in \mathbb{N}, 1 \leq k \leq l \right\}$, where $l$ represents the number of SSL layers, $p$ denotes the preset number of tokens for PT. Therefore, the input and output of the first layer of the Transformer are as follows:

$$[Z_1, E_1] = L_1([P_1, E_0]), \tag{1}$$

where $Z_1 \in \mathbb{R}^{p \times d}$ is the variable generated by the first frozen transformer encoder at the prompt token position, which will be replaced by $P_1 \in \mathbb{R}^{p \times d}$ in the next computation. Thus, the PT calculation for other layers is as follows:

$$[Z_i, E_i] = L_i([P_i, E_{i-1}]), \quad \text{for } i = 2, 3, ..., l. \tag{2}$$

Taking the most commonly used SSL feature in the ADD domain, XLSR-300m[2], as an example, after the final 24 layers, we obtain a matrix output $I = [Z_{24}, E_{24}]$. $I$ will serve as the input to AASIST. The back-end AASIST classifier fully follows the structure of SSL-AASIST by Tak et al. [10], utilizing spectro-temporal graph attention to capture time-frequency features. The final output is a two-dimensional logits score, which is optimized through weighted cross-entropy (WCE) loss.

## 3.2 WPT-SSL-AASIST

In the PT-SSL-AASIST framework, the initial embedding $E_0 \in \mathbb{R}^{t \times d}$, extracted by the SSL front-end, retains high temporal resolution from raw waveform inputs but lacks explicit frequency distribution and cross-type frequency attention ability. To achieve frequency-sensitive modeling for all-type ADD, we proposed WPT-SSL-AASIST, which introduces wavelet prompt tokens to enhance the frequency perception capability of SSL. The difference between WPT and PT lies in the prompt initialization. We use Xavier uniform initialization to initialize two sets of prompt tokens: wavelet initial tokens $\mathbf{T} = \left\{ \mathbf{T}_k \in \mathbb{R}^{w \times d} \mid k \in \mathbb{N}, 1 \leq k \leq l \right\}$ and prompt token $\mathbf{P} = \left\{ \mathbf{P}_k \in \mathbb{R}^{p \times d} \mid k \in \mathbb{N}, 1 \leq k \leq l \right\}$, where $w$ and $p$ denotes the preset number of Wavelet tokens and PT tokens, respectively. For the wavelet initial token, we use the efficient and straightforward wavelet Haar to perform the DWT transformation. Haar wavelets consist of the low-pass filter L, and the high-pass filter H, as follows:

$$L = \frac{1}{\sqrt{2}} [1, 1]^T, H = \frac{1}{\sqrt{2}} [1, -1]^T. \tag{3}$$

We can obtain four sub-bands, which can be expressed as:

$$T_{LL}, \{T_{LH}, T_{HL}, T_{HH}\} = \text{DWT}(T). \tag{4}$$

The Haar wavelet transform generates four components: the low-frequency component (LL), as well as the high frequency in the vertical (LH), horizontal (HL), and diagonal (HH) directions. Each component has a size of $\frac{w}{2} \times \frac{d}{2}$, and then we reshape each component to a size of $\frac{w}{4} \times d$. Based on this operation, each token can correspond to a specific frequency component. Finally, we concatenate LL, LH, HL, and HH components to form the wavelet prompt $\mathbf{W} = \left\{ \mathbf{W}_k \in \mathbb{R}^{w \times d} \mid k \in \mathbb{N}, 1 \leq k \leq l \right\}$. The above process is illustrated in the lower part of Figure 2.

After obtaining the wavelet prompt, we concatenate it with the prompt token $P$ at each layer. Thus, the WPT process can be illustrated as follows:

$$[Z_i, E_i] = L_i([W_i, P_i, E_{i-1}]), \quad \text{for } i = 1, 2, ..., l. \tag{5}$$

Similar to PT-SSL-AASIST, the output of the transformer final layer $I = [Z_l, E_l]$ will be sent to the AASIST backend and trained using the WCE loss.

---

[2] https://huggingface.co/facebook/wav2vec2-xls-r-300m

**Table 1: Statistics of all-type ADD benchmark in terms of training, development, and evaluation set.**

| Type | Source | Train | Dev | Eval |
|---|---|---|---|---|
| Speech | 19LA | 25,380 | 24,844 | 71,237 |
| Sound | Codecfake-A3 | 69,378 | 9,911 | 19,823 |
| Singing | CtrSVDD | 84,404 | 43,625 | 92,769 |
| Music | FakeMusicCaps | 20,861 | 6,058 | 6,122 |
| All | Combined Sources | 199,023 | 84,438 | 189,951 |

## 4 ALL-TYPE ADD BENCHMARK

In this section, we will present the benchmark experimental setup, including the four type ADD datasets used, the CMs employed, the training and testing protocols, and the detailed implementation of the entire experiment.

### 4.1 Dataset

To evaluate CM's ability to detect all types of deepfake audio, the selection of datasets is crucial. The principles for selection include being relatively clean and devoid of partially spoofed scenarios. Our aim is to thoroughly explore the capabilities of CMs in relatively clean environments, as removing other interferences such as noise is beneficial for studying cross-type ADD. Details of the dataset can be found in Table 1.

**Speech-19LA**. The 19LA dataset is one of the most popular datasets in the field of ADD. It contains a total of 12,456 real and 108,978 fake speeches generated using 11 Text-to-Speech (TTS) and 8 Voice Conversion (VC) spoofing algorithms (A01-A19). The real source domain of 19LA from VCTK [42]. The training and development sets consist of data from spoofing systems denoted by A01 to A06, while the evaluation set includes other generation methods denoted by A07 to A19. By testing on data where the generation method does not appear in the training dataset, we can evaluate the effectiveness of the CM.

**Sound-Codecfake-A3**. Considering the principles of data selection and the total volume of the dataset, we chose the Codecfake A3 subset for sound. The real source domain is from the training subset of Audiocaps [43], and the fake sounds are generated using AudioGen based on the corresponding caption. Since some of the original Audiocaps audio links are no longer active, a small number of real sounds could not be downloaded. Ultimately, this condition includes 49,274 real sounds and 49,838 fake sounds. We randomly divided all the sound data into training, validation, and evaluation sets in a ratio of 7:1:2.

**Singing voice-CtrSVDD**. SVDD [4] is the first singing voice detection challenge. The real source domain of CtrSVDD includes multiple open-source Mandarin and Japanese singing datasets, and the synthesis methods encompass 14 different Singing Voice Synthesis (SVS) and Singing Voice Conversion (SVC) methods (A01-A14). The training protocol adheres to the principle that the generation methods used in the evaluation set are not visible in the training set, utilizing A01-A08 for training and A09-A14 for testing. We strictly follow the original training, validation, and evaluation split protocol of CtrSVDD.

**Music-FakeMusicCaps**. FakeMusicCaps [9] is a deepfake music detection dataset. The real source domain of FakeMusicCaps is the MusicCaps [44] dataset, which consists of 5.5k 10-second music

clips from AudioSet [45], each paired with an annotation by a professional musician. The fake-generated captions are identical to those in SunoCaps [46], and the synthesis methods include six approaches, labeled as TTM01-TTM05, along with one unknown method. We have restructured the dataset based on the fake methods. The real music is divided into a 7:1:2 ratio for training, validation, and testing. The training set for fake music includes TTM01-TTM03, the validation set includes TTM04, and the testing set includes TTM05 and the unknown method.

## 4.2 Baseline Countermeasure

In our benchmark, we have five fundamental models, which we refer to as Spec-Resnet, AASIST, MERT-AASIST, WavLM-AASIST, and XLSR-AASIST based on their front-end and back-end concatenation.

Spec-Resnet is a classic audio classification method that combines Spectrogram with Resnet [47]. Although traditional features have previously underperformed compared to SSL features, further research is needed to study their generalizability across different types.

AASIST is a powerful deepfake CM that takes raw waveforms as input. It uses a sinc convolution layer [48] and six residual blocks to extract high-level representations, which are then fed into spectral and temporal graph attention for binary classification.

For SSL-AASIST, we employ three types of SSL features: MERT [3], WavLM [4], and XLSR [5]. MERT is a self-supervised acoustic music understanding model that has achieved state-of-the-art (SOTA) performance in multiple music information retrieval (MIR) tasks. We believe it has the potential to distinguish deepfake music. When WavLM was proposed, it exhibited the SOTA performance among SSL models for all speech downstream tasks. Its performance in the ADD domain, including all-type ADD task, warrants further investigation. XLSR is the most representative of SSL for ADD tasks, and it has been proven to be the best SSL feature in the ADD domain across multiple SSL surveys [19, 49].

Regarding the training paradigms for SSL-AASIST, this paper explores four different approaches: FR-SSL-AASIST, FT-SSL-AASIST, PT-SSL-AASIST, and WPT-SSL-AASIST. Among these, FR and FT represent two distinct training paradigms: freezing and fine-tuning. In the FR method, all SSL parameters remain frozen, with only the final-layer SSL features fed into the AASIST back-end. FT method follows the same approach as FR but allows the SSL parameters to be updated through gradient backpropagation during training. The PT and WPT approach we have introduced in Section 3.

## 4.3 Training and Evaluation Protocol

To evaluate the all-type ADD capability of CM, we first conducted single-type training experiments, where the model was trained on one type of ADD dataset and tested on other types. In these experiments, the five CMs mentioned in the previous section were trained using a single type of training set and tested separately on each type of test set. For SSL-AASIST, different training paradigms can be employed, including FR, FT, PT, and WPT. To further address the all-type ADD task, we conducted all-type co-training experiments.

[3]https://huggingface.co/m-a-p/MERT-v1-330M
[4]https://huggingface.co/microsoft/wavlm-large
[5]https://huggingface.co/facebook/wav2vec2-xls-r-300m

**Table 2: EER (%) results of the countermeasures (frozen SSL) trained on single-type ADD training set.**

| Train | Countermeasure | Speech | Sound | Singing | Music | AVG |
|---|---|---|---|---|---|---|
| Speech | Spec-Resnet | 5.58 | 48.64 | 45.15 | 47.01 | 36.60 |
| Speech | AASIST | 1.48 | 48.32 | 40.71 | 47.75 | 34.57 |
| Speech | FR-MERT-AASIST | 4.80 | **47.60** | 44.51 | 48.89 | 36.45 |
| Speech | FR-WavLM-AASIST | 2.49 | 47.96 | 38.67 | **42.75** | 32.97 |
| Speech | FR-XLSR-AASIST | **1.28** | 49.51 | **29.72** | 49.82 | **32.58** |
| Sound | Spec-Resnet | 49.67 | 8.87 | 47.77 | 44.22 | 37.63 |
| Sound | AASIST | 37.39 | **0.43** | 42.56 | **10.44** | 22.71 |
| Sound | FR-MERT-AASIST | 23.37 | 0.64 | 43.31 | 49.82 | 29.29 |
| Sound | FR-WavLM-AASIST | 39.25 | 7.09 | 36.67 | 46.47 | 32.37 |
| Sound | FR-XLSR-AASIST | **16.88** | 2.40 | **31.82** | 33.65 | **21.19** |
| Singing | Spec-Resnet | 37.54 | 46.04 | 23.59 | **32.70** | 34.97 |
| Singing | AASIST | 33.06 | 38.23 | 20.51 | 36.62 | 32.11 |
| Singing | FR-MERT-AASIST | 43.86 | 42.88 | 29.95 | 44.24 | 40.23 |
| Singing | FR-WavLM-AASIST | 16.19 | 41.80 | 18.74 | 39.18 | 28.98 |
| Singing | FR-XLSR-AASIST | **12.89** | **34.41** | **9.45** | 35.87 | **23.16** |
| Music | Spec-Resnet | 46.33 | 47.52 | 48.33 | 15.61 | 39.45 |
| Music | AASIST | 31.81 | 47.26 | 44.12 | 8.36 | 32.89 |
| Music | FR-MERT-AASIST | **27.88** | 44.45 | **34.56** | 7.62 | **28.63** |
| Music | FR-WavLM-AASIST | 45.88 | 43.64 | 45.15 | 15.80 | 37.62 |
| Music | FR-XLSR-AASIST | 48.89 | **40.54** | 43.41 | 9.67 | 35.63 |

Specifically, we trained the CM using all types of training set and tested it on each type of evaluation set.

## 4.4 Implementation Details

For the pre-processing of the ADD baseline models, all audio samples were first down-sampled to 16,000 Hz and trimmed or padded to 64600 samples (same as the original AASIST and SSL-AASIST). For the Spec-Resnet, the spectrogram was computed with the number of FFT points set to 512, the hop length set to 160, and the window length set to 512. The back-end Resnet used Resnet18 followed by a fully connected layer to down-sample to 2 dimensions. For the training paradigm, FT-SSL-AASIST adopted the training parameters from Tak et al. [10], with an initial learning rate of $10^{-6}$ and a batch size of 14. For FR, PT, and WPT, an initial learning rate of $5^{-4}$ and a batch size of 32 were used. The dimensions of the 4-second audio processed by the SSL feature extractor are (201, 1024). For PT, the optimal number of prompt tokens was 10, while for WPT, the number of wavelet tokens was 4 and the number of regular prompt tokens was 6. In the single-type training experiments, the epoch was set to 50, and the learning rate was halved every 10 steps. In the co-training experiments, the epoch was set to 20, and the learning rate was halved every 4 steps.

## 5 EXPERIMENTS

### 5.1 Investigation for Single-Type Training

In this section, we used a single-type training set to train deepfake CMs. First, we validated the performance of Spec-Resnet, AASIST, and FR-SSL-AASIST, as shown in Table 2. By analyzing the frozen

**Table 3: EER (%) results of the countermeasures (finetuned SSL) trained on single-type ADD training set.**

| Train | Countermeasure | Speech | Sound | Singing | Music | AVG |
|---|---|---|---|---|---|---|
| Speech | FT-Mert-AASIST | 6.99 | 48.37 | 48.86 | 44.43 | 37.16 |
| Speech | FT-WavLM-AASIST | 1.50 | **44.62** | 35.77 | 42.19 | 31.02 |
| Speech | FT-XLSR-AASIST | **0.38** | 49.57 | **29.76** | **31.01** | **27.68** |
| Sound | FT-MERT-AASIST | 21.69 | **0.20** | 48.23 | 43.68 | 28.45 |
| Sound | FT-WavLM-AASIST | 32.10 | **0.20** | 47.76 | **21.72** | 25.45 |
| Sound | FT-XLSR-AASIST | **9.22** | 0.21 | **35.96** | 44.77 | **22.54** |
| Singing | FT-MERT-AASIST | 43.51 | 41.92 | 30.58 | 41.81 | 39.46 |
| Singing | FT-WavLM-AASIST | 13.07 | 36.68 | 8.00 | 40.32 | 24.52 |
| Singing | FT-XLSR-AASIST | **7.56** | **31.08** | **5.60** | 37.36 | **20.40** |
| Music | FT-MERT-AASIST | **24.03** | 46.50 | 44.79 | **15.53** | **32.21** |
| Music | FT-WavLM-AASIST | 48.82 | 48.82 | 46.22 | 47.03 | 47.72 |
| Music | FT-XLSR-AASIST | 39.03 | 47.99 | 47.93 | 48.70 | 45.91 |

**Table 4: EER (%) comparison with different number of token.**

| Token | Parm | Speech | Sound | Singing | Music | AVG |
|---|---|---|---|---|---|---|
| 2 | 0.50M | 0.75 | 45.29 | 35.00 | 42.71 | 30.94 |
| 10 | 0.69M | **0.22** | 47.26 | **33.84** | 41.85 | **30.79** |
| 20 | 0.94M | 0.58 | 44.11 | 43.35 | 41.64 | 32.42 |
| 100 | 2.90M | 3.01 | **37.05** | 49.41 | **35.66** | 31.28 |
| 200 | 5.36M | 4.99 | 44.45 | 47.61 | 36.37 | 33.36 |

**Table 5: EER (%) comparison with different paradigms.**

| Paradigm | Speech | Sound | Singing | Music | AVG |
|---|---|---|---|---|---|
| Shallow-PT | 0.75 | **45.29** | 39.87 | 44.24 | 32.54 |
| After-PT | 0.53 | 46.88 | 41.55 | 44.05 | 33.25 |
| Del-PT | 0.72 | 47.23 | 41.45 | 42.87 | 33.07 |
| PT | **0.22** | 47.26 | **33.84** | **41.85** | **30.79** |



**Figure 3: Different paradigms of PT-SSL-AASIST.**

SSL features, we can observe the inherent deepfake detection capabilities of SSL. Firstly, for the speech-trained CM, the best performance was achieved by XLSR-AASIST, which obtained the lowest in-domain EER of 1.28% and the lowest average EER of 32.58%. This result aligns with our expectations. Additionally, there is a significant drop in performance for cross-type ADD task. We observed that the singing test resulted in an EER of 29.72%, which is significantly better than the nearly 50% EER for sound and music. This denotes that there may be some shared characteristics in distinguishing between deepfake speech and deepfake singing. For the Sound task, AASIST achieved the lowest ID EER of 0.43%. However, its poor performance on speech and singing resulted in an average performance that was inferior to XLSR-AASIST. We believe that AASIST's strong performance on sound is due to the fact that Codecfake-A3 uses only one deepfake method, leading to some degree of over-fitting. Nonetheless, its performance on music suggests that this over-fitting also benefits music detection. This indicates that there may be some commonality in distinguishing between sound and music, which are types of audio that usually do not contain human voices. For the singing task, XLSR-AASIST achieved both the best ID and best average performance. It obtained an ID EER of 9.45% and an average EER of 23.16%. Notably, XLSR-AASIST, which was trained solely on singing data, also demonstrated an EER of 12.89% on the speech task. This further indicates the commonality between speech and singing. For the music task, MERT-AASIST achieved the best ID and average performance with EER of 7.62% and 28.63%, respectively. This aligns with our expectations upon introducing MERT features in ADD task.
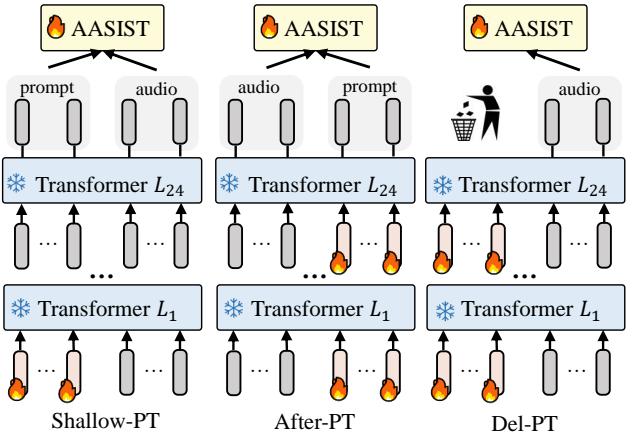
Then, we investigate the SSL-AASIST through fine-tuning full SSL layer as shown in Table 3. Overall, the final results were completely consistent with the frozen SSL models. For the speech-trained, sound-trained, and singing-trained CMs, the best performance was achieved by FT-XLSR-AASIST, with average EERs of 27.68%, 22.54%, and 20.40%, respectively. For the music-trained CMs, the best performance was achieved by FT-MERT-AASIST, with an average EER of 32.21%. It is also noteworthy that FT-SSL-AASIST consistently achieved lower EERs across various ID tasks compared to FR-SSL-AASIST. For instance, FT-XLSR-AASIST achieved EER
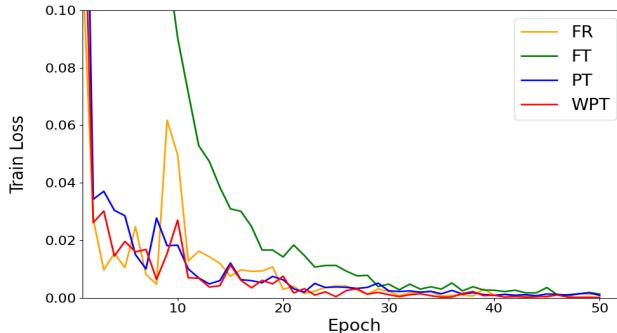
of 0.38% for speech, 0.21% for sound, and 5.60% for singing, representing reductions of 0.9%, 2.19%, and 3.85% respectively compared to FR-SSL-AASIST. However, for music-trained SSL, all features exhibited a decline compared to FR, highlighting the challenges of fine-tuning. This indicates that fine-tuning, while requiring extensive parameter training, may also necessitate setting different hyper-parameters based on the type of data and SSL.

## 5.2 Prompt Tuning Countermeasures

**PT-SSL-AASIST.** To evaluate the effectiveness of PT and determine their optimal parameters, we integrated PT into the XLSR-AASIST, which performed best in the previous section, trained on the speech dataset. There are two aspects worth investigating for PT: the preset number of tokens for PT and the paradigm for PT (connection method, prompt position, etc.). We conducted ablation experiments on the number of tokens and the paradigm, as shown in the Table 4 and Table 5, respectively.

**Table 6: EER (%) and training parameters comparison with different paradigms of speech-trained XLSR-AASIST.**

| Countermeasure | Parm | Speech | Sound | Singing | Music | AVG |
|---|---|---|---|---|---|---|
| FR-XLSR-AASIST | 0.45M | 1.28 | 49.51 | **29.72** | 49.82 | 32.58 |
| FT-XLSR-AASIST | 315.89M | 0.38 | 49.57 | 29.76 | 31.01 | 27.68 |
| PT-XLSR-AASIST | 0.69M | 0.22 | 47.26 | 33.84 | 41.85 | 30.79 |
| WPT-XLSR-AASIST | 0.69M | **0.15** | **45.36** | 33.32 | **28.61** | **26.86** |



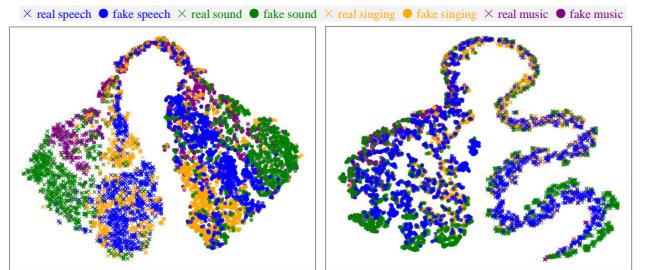**Figure 4: Convergence speed of different training paradigms.**

For the number of tokens in PT, we experimented with 2, 10, 20, 100, and 200. The results showed that when the token number was set to 10, the best speech test set EER of 0.22% and the lowest average EER of 30.79% were achieved. As the number of tokens increased, the parameter count for PT training also increased, but the effectiveness decreased. This is may due to the fact that the number of audio tokens is 201, and an excessive number of prompt tokens can cause the audio tokens to become sparse, hindering the learning of the audio's inherent information.

After determining the number of tokens, we investigated three paradigms for PT-SSL-AASIST, including Shallow-PT, After-PT, and Del-PT, as shown in Figure 3. Shallow-PT refers to inserting learnable prompts only in the first transformer encoder layer, which can demonstrate the importance of the deep paradigm where prompts are inserted in each layer. After-PT places the prompt position after the audio token, which might be effective due to the artifact information located in the silent region at the beginning of the audio [50]. Del-PT involves deleting the prompt token in the last layer, using only the audio tokens for classification, a method considered effective in some vision tasks [37]. Experimental results indicate that our proposed PT-SSL-AASIST paradigm is optimal, where prompts are inserted in each layer and the final layer combines the prompt and audio tokens for input into AASIST.

**WPT-SSL-AASIST.** After deciding the PT architecture, we introduced WPT, applying DWT to the first four of the ten prompt tokens to better capture the frequency information of the audio. We compared the performance of FR, FT, PT, and WPT using the speech-trained XLSR-AASIST, as shown in Table 6. It can be observed that WPT achieved the best results compared to PT, obtaining a 0.15% EER on the ID speech evaluation set, with an average EER of 27.55%.

**Table 7: EER (%) results for the countermeasures co-trained on the complete ADD training set.**
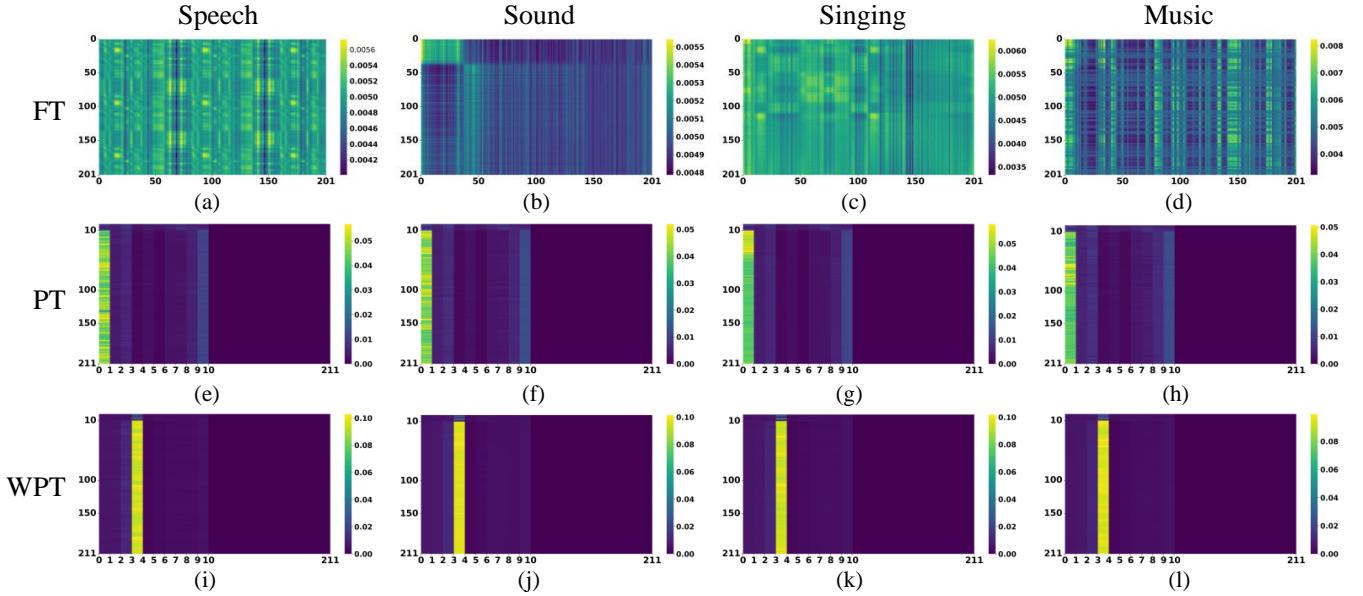
| Countermeasure | Speech | Sound | Singing | Music | AVG |
|---|---|---|---|---|---|
| Spec-Resnet | 29.37 | 23.37 | 37.17 | 42.75 | 33.17 |
| AASIST | 3.78 | 0.86 | 20.01 | 11.70 | 9.09 |
| FR-WavLM-AASIST | 3.44 | 10.21 | 17.83 | 26.02 | 14.38 |
| FT-WavLM-AASIST | **1.31** | 2.53 | 16.48 | 22.90 | 10.81 |
| PT-WavLM-AASIST | 3.09 | 8.81 | 15.84 | **16.73** | 11.12 |
| WPT-WavLM-AASIST | 2.04 | **1.10** | 9.28 | 18.21 | **7.66** |
| FR-MERT-AASIST | **2.90** | 4.60 | **12.14** | 24.91 | 11.14 |
| FT-MERT-AASIST | 6.24 | 1.17 | 31.67 | 13.77 | 13.21 |
| PT-MERT-AASIST | 6.06 | 1.28 | 32.59 | 9.29 | 12.31 |
| WPT-MERT-AASIST | 6.59 | **1.01** | 22.68 | **8.53** | **9.70** |
| FR-XLSR-AASIST | 3.02 | 5.45 | 10.86 | 22.67 | 10.50 |
| FT-XLSR-AASIST | 1.77 | **0.49** | 8.93 | 8.71 | 4.98 |
| PT-XLSR-AASIST | 2.00 | 1.11 | 14.54 | 9.29 | 6.74 |
| WPT-XLSR-AASIST | **0.72** | 1.29 | **7.47** | **4.83** | **3.58** |



**Figure 5: T-SNE visualization for FT-XLSR-AASIST (left) and WPT-XLSR-AASIST (right). Different colors indicate features from different types: blue=speech, green=sound, orange=singing, purple=music. Different shapes represent different categories: cross=real, point=fake.**

Moreover, WPT does not increase the number of training parameters compared to PT, and compared to FT, the training parameters are reduced by 458 times. Overall, WPT outperformed FT, which in turn outperformed PT, and PT outperformed FR. We also recorded the training convergence speeds, as shown in Figure 4. It can be observed that FR, PT, and WPT converged significantly faster than FT, and both FT and PT showed less fluctuation compared to FR during convergence.

### 5.3 Co-trained Countermeasures

Although the speech-trained WPT-XLSR-AASIST achieved extremely low EER on speech deepfake test set, it still exhibited significant performance degradation on detecting deepfake sound, singing, and music. Therefore, we began to investigate co-trained CM, combining the training sets of the four types to achieve all-type audio deepfake CM.

**Figure 6: Attention map of the final transformer in co-trained XLSR-AASIST. Each column corresponds to the same deepfake audio sample. For both PT and WPT, we magnified the position of prompt tokens (1-10).**

The results of the co-training experiment are shown in Table 7. Firstly, the effectiveness of the data-driven approach can be observed, with a significant reduction in average EER compared to single-type trained CMs. The best performing SSL in the co-training experiment is the XLSR-AASIST. For the XLSR-AASIST training paradigm, WPT outperformed FT, PT, and FR, achieving EERs of 3.58%, 4.98%, 6.74%, and 10.50%, respectively. This training paradigm's performance aligns with that of the speech-trained XLSR-AASIST shown in Table 6. Notably, WPT consistently achieves the best performance across different SSL features. For instance, WPT-WavLM-AASIST and WPT-MERT-AASIST achieve EER of 7.66% and 9.70%, respectively.

## 5.4 Interpretability

**Type Invariance in T-SNE Visualization.** To further understand the interpretability of the WPT training paradigm, we first performed T-SNE visualization on the embeddings before the final fully connected layer of AASIST. Specifically, we applied T-SNE visualization to the embeddings from the co-trained FT-XLSR-AASIST and WPT-XLSR-AASIST on evaluation sets of four audio types. For each type, we selected 2,000 samples randomly, comprising 1,000 genuine samples and 1,000 fake samples. The results are presented in Figure 5. Firstly, it can be observed that both FT and WPT are capable of separating the test real and fake samples. However, there is a notable difference. FT demonstrates distinct clustering within both the genuine and fake regions, where speech, sound, singing, and music samples form separate clusters. In contrast, WPT does not exhibit such separation within either the genuine or fake regions, resulting in overlap among the four types. This indicates that WPT maintains type invariance when performing the all-type ADD task.

**Type Invariance in Attention Distribution.** To further investigate the intrinsic differences in training paradigms for detecting deepfakes, we plotted the attention maps of the final transformer encoder layer, as shown in Figure 6. It is evident that FT exhibits different attention distributions when processing different types of audio. Interestingly, the attention patterns for speech and singing are similar, exhibiting overall high values with some regions of exceptionally high intensity. The attention patterns for sound and music are also similar, displaying a mix of high and low values in all region. This observation is consistent with the experimental results from single-type training. For the PT and WPT paradigms, we can observe consistency in their detecting of different types. The PT paradigm focuses on the first prompt token, but the values are not high, and there are noticeable value changes when dealing with different types, with some attention also present on the 10th prompt token. In contrast, WPT paradigm demonstrates significant invariance in detecting diverse audio types, with a focus on the 4th token corresponding to the wavelet HH token, which determines high-frequency details through diagonal orientation analysis.

## 5.5 Conclusion

In this paper, we are dedicated to studying the all-type ADD task. We are the first to establish a comprehensive benchmark for evaluating the performance of current CMs on the all-type ADD task. Building on this foundation, we propose the PT-SSL training paradigm, which maintains performance while significantly reducing training parameters. Finally, to achieve all-type CM, we propose the WPT-SSL training paradigm, which leverages wavelet prompts to capture the type-invariant auditory deepfake information of SSL features. Our proposed co-trained WPT-XLSR-AASIST achieves an average EER of 3.58% across all-type ADD evaluation set.

# REFERENCES

[1] Massimiliano Todisco, Xin Wang, Ville Vestman, Md. Sahidullah, Héctor Delgado, Andreas Nautsch, Junichi Yamagishi, Nicholas Evans, Tomi H. Kinnunen, and Kong Aik Lee. ASVspoof 2019: Future Horizons in Spoofed and Fake Audio Detection. In *Proc. Interspeech 2019*, pages 1008–1012, 2019.

[2] Xuechen Liu, Xin Wang, Md Sahidullah, Jose Patino, Héctor Delgado, Tomi Kinnunen, Massimiliano Todisco, Junichi Yamagishi, Nicholas Evans, Andreas Nautsch, et al. Asvspoof 2021: Towards spoofed and deepfake speech detection in the wild. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2023.

[3] Xin Wang, Héctor Delgado, Hemlata Tak, Jee-weon Jung, Hye-jin Shim, Massimiliano Todisco, Ivan Kukanov, Xuechen Liu, Md Sahidullah, Tomi Kinnunen, Nicholas Evans, Kong Aik Lee, and Junichi Yamagishi. ASVspoof 5: Crowdsourced speech data, deepfakes, and adversarial attacks at scale. In *ASVspoof Workshop 2024 (accepted)*, 2024.

[4] You Zhang, Yongyi Zang, Jiatong Shi, Ryuichi Yamamoto, Tomoki Toda, and Zhiyao Duan. Svdd 2024: The inaugural singing voice deepfake detection challenge. In *2024 IEEE Spoken Language Technology Workshop (SLT)*, pages 782–787. IEEE, 2024.

[5] Yuankun Xie, Jingjing Zhou, Xiaolin Lu, Zhenghao Jiang, Yuxin Yang, Haonan Cheng, and Long Ye. Fsd: An initial chinese dataset for fake song detection. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4605–4609. IEEE, 2024.

[6] Yongyi Zang, Jiatong Shi, You Zhang, Ryuichi Yamamoto, Jionghao Han, Yuxun Tang, Shengyuan Xu, Wenxiao Zhao, Jing Guo, Tomoki Toda, et al. Ctrsvdd: A benchmark dataset and baseline analysis for controlled singing voice deepfake detection. *arXiv preprint arXiv:2406.02438*, 2024.

[7] Zeyu Xie, Baihan Li, Xuenan Xu, Zheng Liang, Kai Yu, and Mengyue Wu. Fakesound: Deepfake general audio detection. In *Proc. Interspeech 2024*, pages 112–116, 2024.

[8] Yuankun Xie, Yi Lu, Ruibo Fu, Zhengqi Wen, Zhiyong Wang, Jianhua Tao, Xin Qi, Xiaopeng Wang, Yukun Liu, Haonan Cheng, et al. The codecfake dataset and countermeasures for the universally detection of deepfake audio. *IEEE Transactions on Audio, Speech and Language Processing*, 2025.

[9] Luca Comanducci, Paolo Bestagini, and Stefano Tubaro. Fakemusiccaps: a dataset for detection and attribution of synthetic music generated via text-to-music models. *arXiv preprint arXiv:2409.10684*, 2024.

[10] Hemlata Tak, Massimiliano Todisco, Xin Wang, Jee-weon Jung, Junichi Yamagishi, and Nicholas Evans. Automatic speaker verification spoofing and deepfake detection using wav2vec 2.0 and data augmentation. In *The Speaker and Language Recognition Workshop (Odyssey 2022)*. ISCA, 2022.

[11] Arun Babu, Changhan Wang, Andros Tjandra, Kushal Lakhotia, Qiantong Xu, Naman Goyal, Kritika Singh, Patrick von Platen, Yatharth Saraf, Juan Pino, Alexei Baevski, Alexis Conneau, and Michael Auli. Xls-r: Self-supervised cross-lingual speech representation learning at scale. In *Interspeech 2022*, pages 2278–2282, 2022.

[12] Mahyar Gohari, Davide Salvi, Paolo Bestagini, and Nicola Adami. Audio features investigation for singing voice deepfake detection. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2025.

[13] Zhiyong Wang, Ruibo Fu, Zhengqi Wen, Jianhua Tao, Xiaopeng Wang, Yuankun Xie, Xin Qi, Shuchen Shi, Yi Lu, Yukun Liu, et al. Mixture of experts fusion for fake audio detection using frozen wav2vec 2.0. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2025.

[14] Jee-weon Jung, Hee-Soo Heo, Hemlata Tak, Hye-jin Shim, Joon Son Chung, Bong-Jin Lee, Ha-Jin Yu, and Nicholas Evans. Aasist: Audio anti-spoofing using integrated spectro-temporal graph attention networks. In *Proceedings of the ICASSP*, pages 6367–6371, 2022.

[15] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *European conference on computer vision*, pages 709–727. Springer, 2022.

[16] Sam Norman-Haignere, Nancy G Kanwisher, and Josh H McDermott. Distinct cortical pathways for music and speech revealed by hypothesis-free voxel decomposition. *neuron*, 88(6):1281–1296, 2015.

[17] Alexander JE Kell, Daniel LK Yamins, Erica N Shook, Sam V Norman-Haignere, and Josh H McDermott. A task-optimized neural network replicates human auditory behavior, predicts brain responses, and reveals a cortical processing hierarchy. *Neuron*, 98(3):630–644, 2018.

[18] Rungsun Munkong and Biing-Hwang Juang. Auditory perception and cognition. *IEEE signal processing magazine*, 25(3):98–117, 2008.

[19] Orchid Chetia Phukan, Gautam Siddharth Kashyap, Arun Balaji Buduru, and Rajesh Sharma. Heterogeneity over homogeneity: Investigating multilingual speech pre-trained models for detecting audio deepfake. *arXiv preprint arXiv:2404.00809*, 2024.

[20] Yassine El Kheir, Youness Samih, Suraj Maharjan, Tim Polzehl, and Sebastian Möller. Comprehensive layer-wise analysis of ssl models for audio deepfake

[21] Qishan Zhang, Shuangbing Wen, and Tao Hu. Audio deepfake detection with self-supervised xls-r and sls classifier. In *ACM Multimedia 2024*, 2024.

[22] Zihan Pan, Tianchi Liu, Hardik B Sailor, and Qiongqiong Wang. Attentive merging of hidden embeddings from pre-trained speech model for anti-spoofing detection. In *Proc. Interspeech 2024*, pages 2090–2094, 2024.

[23] Zirui Zhang, Wei Hao, Aroon Sankoh, William Lin, Emanuel Mendiola-Ortiz, Junfeng Yang, and Chengzhi Mao. I can hear you: Selective robust training for deepfake audio detection. In *The Thirteenth International Conference on Learning Representations*, 2025.

[24] Piotr Kawa, Marcin Plata, and Piotr Syga. Defense against adversarial attacks on audio deepfake detection. In *Proc. Interspeech 2023*, pages 5276–5280, 2023.

[25] Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation. In *The Eleventh International Conference on Learning Representations*, 2022.

[26] Anmol Guragain, Tianchi Liu, Zihan Pan, Hardik B. Sailor, and Qiongqiong Wang. Speech foundation model ensembles for the controlled singing voice deepfake detection (ctrsvdd) challenge 2024. In *2024 IEEE Spoken Language Technology Workshop (SLT)*, pages 774–781, 2024.

[27] Qishan Zhang, Shuangbing Wen, Fangke Yan, Tao Hu, and Jun Li. Xwsb: A blend system utilizing xls-r and wavlm with sls classifier detection system for svdd 2024 challenge. In *2024 IEEE Spoken Language Technology Workshop (SLT)*, pages 788–794. IEEE, 2024.

[28] Xuanjun Chen, Haibin Wu, Roger Jang, and Hung-yi Lee. Singing voice graph modeling for singfake detection. In *Proc. Interspeech 2024*, pages 4843–4847, 2024.

[29] Yizhi Li, Ruibin Yuan, Ge Zhang, Yinghao Ma, Xingran Chen, Hanzhi Yin, Chenghao Xiao, Chenghua Lin, Anton Ragni, Emmanouil Benetos, et al. Mert: Acoustic music understanding model with large-scale self-supervised training. In *ICLR*, 2024.

[30] Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1505–1518, 2022.

[31] Yupei Li, Qiyang Sun, Hanqian Li, Lucia Specia, and Björn W Schuller. Detecting machine-generated music with explainability–a challenge and early benchmarks. *arXiv preprint arXiv:2412.13421*, 2024.

[32] Zhaolin Wei, Dengpan Ye, Jiacheng Deng, and Yuhan Lin. From voices to beats: Enhancing music deepfake detection by identifying forgeries in background. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2025.

[33] Yupei Li, Manuel Milling, Lucia Specia, and Björn W Schuller. From audio deepfake detection to ai-generated music detection–a pathway and overview. *arXiv preprint arXiv:2412.00571*, 2024.

[34] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, 2021.

[35] Xiaohua Zhai, Joan Puigcerver, Alexander Kolesnikov, Pierre Ruyssen, Carlos Riquelme, Mario Lucic, Josip Djolonga, Andre Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, et al. A large-scale study of representation learning with the visual task adaptation benchmark. *arXiv preprint arXiv:1910.04867*, 2019.

[36] Seungryong Yoo, Eunji Kim, Dahuin Jung, Jungbeom Lee, and Sungroh Yoon. Improving visual prompt tuning for self-supervised vision transformers. In *International Conference on Machine Learning*, pages 40075–40092. PMLR, 2023.

[37] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. In *The Twelfth International Conference on Learning Representations*.

[38] Sheng Shen, Shijia Yang, Tianjun Zhang, Bohan Zhai, Joseph E. Gonzalez, Kurt Keutzer, and Trevor Darrell. Multitask vision-language prompt tuning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 5656–5667, January 2024.

[39] Kai-Wei Chang, Haibin Wu, Yu-Kai Wang, Yuan-Kuei Wu, Hua Shen, Wei-Cheng Tseng, Iu-thing Kang, Shang-Wen Li, and Hung-yi Lee. Speechprompt: Prompting speech language models for speech processing tasks. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.

[40] Yuzhuo Liu, Xubo Liu, Yan Zhao, Yuanyuan Wang, Rui Xia, Pingchuan Tain, and Yuxuan Wang. Audio prompt tuning for universal sound separation. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1446–1450, 2024.

[41] Hideyuki Oiso, Yuto Matsunaga, Kazuya Kakizaki, and Taiki Miyagawa. Prompt tuning for audio deepfake detection: Computationally efficient test-time domain adaptation with limited target dataset. In *Interspeech 2024*, pages 2710–2714, 2024.

[42] Christophe Veaux, Junichi Yamagishi, and Kirsten MacDonald. Cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit. 2017.

[43] Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. Audiocaps: Generating captions for audios in the wild. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational

*Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 119–132, 2019.

[44] Andrea Agostinelli, Timo I Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, et al. Musiclm: Generating music from text. *arXiv preprint arXiv:2301.11325*, 2023.

[45] Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 776–780. IEEE, 2017.

[46] M Civit, V Drai-Zerbib, D Lizcano, and MJ Escalona. Sunocaps: A novel dataset of text-prompt based ai-generated music with emotion annotations. *Data in Brief*, 55:110743–110743, 2024.

[47] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[48] Mirco Ravanelli and Yoshua Bengio. Speaker recognition from raw waveform with sincnet. In *2018 IEEE spoken language technology workshop (SLT)*, pages 1021–1028. IEEE, 2018.

[49] Octavian Pascu, Adriana Stan, Dan Oneata, Elisabeta Oneata, and Horia Cucu. Towards generalisable and calibrated audio deepfake detection with self-supervised representations. In *Interspeech 2024*, pages 4828–4832, 2024.

[50] Yuxiang Zhang, Zhuo Li, Jingze Lu, Hua Hua, Wenchao Wang, and Pengyuan Zhang. The impact of silence on speech anti-spoofing. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2023.