# Voice Spoofing Countermeasures: Taxonomy, State-of-the-art, experimental analysis of generalizability, open challenges, and the way forward

Awais Khan[1], Khalid Mahmood Malik[1], James Ryan[1], and Mikul Saravanan[1]

[1]Department of Computer Science and Engineering, Oakland University, Rochester, MI 48309, USA

arXiv:2210.00417v2 [eess.AS] 21 Nov 2022

*Abstract*—**Malicious actors may seek to use different voice spoofing attacks to fool ASV systems and even use them for spreading misinformation. Various countermeasures have been proposed to detect these spoofing attacks. Due to the extensive work done on spoofing detection in automated speaker verification (ASV) systems in the last 6-7 years, there is a need to classify the research and perform qualitative and quantitative comparisons on state-of-the-art countermeasures. Additionally, no existing survey paper has reviewed integrated solutions to voice spoofing evaluation and speaker verification, adversarial/anti-forensics attacks on spoofing countermeasures and ASV itself, or unified solutions to detect multiple attacks using a single model. Further, no work has been done to provide an apples-to-apples comparison of published countermeasures in order to assess their generalizability by evaluating them across corpora. In this work, we conduct a review of the literature on spoofing detection using hand-crafted features, deep learning, end-to-end, and universal spoofing countermeasure solutions to detect speech synthesis (SS), voice conversion (VC), and replay attacks. Additionally, we also review integrated solutions to voice spoofing evaluation and speaker verification, adversarial and anti-forensics attacks on voice countermeasures, and ASV. The limitations and challenges of the existing spoofing countermeasures are also presented. We report the performance of these countermeasures on several datasets and evaluate them across corpora. For the experiments, we employ the ASVspoof2019 and VSDC datasets along with GMM, SVM, CNN, and CNN-GRU classifiers. (For reproduce-ability of the results, the code of the testbed can be found at our GitHub Repository[1]).**

*Index Terms*—**Voice Spoofing System, Anti Spoofing,Voice Spoofing Countermeasures, Comparative Analysis, Speech Spoofing Detection.**

## I. Introduction

AUTOMATIC speaker verification systems (ASVs) are now in wide use for authentication, such as for over-the-phone banking, and are becoming an increasingly important biometric solution in the current climate of the pandemic to limit the spread of disease. Although ASVs can be used to validate an identity, they are prone to spoofing attacks, which may allow an unauthorized user to gain access to privileged information such as a bank account. Some of these attacks, such as speech synthesis and voice conversion, are also being used in cyberspace to spread misinformation and

disinformation. In addition, ASV systems are more prevalent in IoT devices, such as smart speakers. A spoofing attack on IoT devices connected in a chain to control someone's home or other devices operated via the ASV system has the potential to allow unauthorized access to security devices, e.g., IoT door locks [54].

Spoofing attacks on ASVs can be grouped into Physical Access (PA) attacks and Logical Access (LA) attacks [4]. An example of a PA attack is a replay attack, where the hacker records the original speaker's voice using any recording device without consent. Later, this prerecorded voice can be replayed onto the ASV system in order to compromise the voice bio-metric based access control in financial institutions and smart homes [54]. LA attacks, conversely, are comprised of artificial, or machine-generated, cloned samples. LA attacks consist of text-to-speech (TTS) synthesis and voice conversion (VC). For TTS synthesis, cloned samples are generated using original speech and transcripts, while VC attacks use only voice samples of the original speaker to train different deep learning models, i.e., a vocoder, waveform, etc., to generate the cloned samples. Again, the objective is to generate realistic voice samples of a target speaker in order to compromise the security of an ASV system, and thus gain access to someone's home, bank account, or other voice-controlled application. Due to a massive escalation in the availability of high-quality recording equipment and the ease of their creation, replay attacks are being generated more easily and rapidly, even by less tech-savvy people, and can be classified into single-hop or multi-hop attacks [6]. More specifically, a single-hop replay attack utilizes a single recording device to play the audio sample of a verified user. In contrast, a multi-hop replay attack utilizes multiple recording devices (i.e., a verified user's voice is recorded, then replayed to another recording device, and then replayed to the ASV system) in a chained scenario. In our prior work [28], we reported a new hybrid voice spoofing attack, i.e., a cloned replay, where the voice of a target speaker is captured and replayed to an ASV system, that can also be generated to spoof an ASV system. A more recent threat to ASVs and spoofing counter measures is adversarial machine learning [75], where machine learning algorithms may be deceived by injecting minor distortions into the audio sample. According to Szegedy et al. [75], deep learning model predictions can be

---

easily modified by extremely small perturbations.

Several countermeasures have been proposed to defeat voice spoofing attacks, and these are often comprised of two parts: the first one (front end) is the feature representation scheme for the input speech signal, and the second one (back end) is a classifier to distinguish between bonafide and spoofed samples. The feature descriptor (front end) should be capable of effectively capturing the traits of the dynamic vocal tracts of a bonafide speaker. Similarly, the back-end classifier should be able to better learn the distinct traits of bonafide and spoofed speech samples in order to accurately discriminate against the spoofed speech samples. In contrast to the traditional approach of front-end and back-end solutions, in the past few years, the research community has focused on deep learning and end-to-end solutions to combat voice spoofing attacks. Existing countermeasures have only tackled single-hop spoofing attacks, and chained-replay attacks and their countermeasures have been largely overlooked [55]. Recent efforts toward in finding a unified solution to the problem of speaker verification, where a single countermeasure may be applied to several attacks, have been studied. Although these countermeasure solutions are only beginning to be explored, there is a marked need for a unified solution as the way forward in ASV anti-spoofing techniques [111] [35]. The surveys on ASV systems and spoof detection techniques conducted to date have primarily been focused on specific spoofing attacks and the countermeasures for them. To the best of our knowledge, six articles [99], [64], [69], [36], [79], [56] have been presented as surveys, and two studies [70] and [19] conducted a comparative analysis of voice spoofing countermeasures.

In a survey of voice spoofing countermeasures, Wu et al. (2015), in [99], [101], [102], provided a good classification of voice spoofing attacks known at the time, along with ASV system vulnerabilities. However, the age of this review means the advanced spoofing attacks that have grown prevalent in ASV systems were not addressed. In addition, is focused on a dedicated countermeasure for each sort of attack, rather than exploring a more unified solution. Following this, Sahidullah et al., [70] and Font et al., [19] published comparative analyses of the spoofing countermeasures that focused solely on replay attacks. According to their work, published during 2017–2018, replay attacks attracted the most attention in the wake of the ASVspoof 2017 challenge. The results of [70], using the ASVspoof2015 dataset, demonstrated that an ensemble of acoustic features, together with machine learning models, produced more accurate results than individual classifiers. Kamble et al. [36] also published a work that covered a subset of the specific speech corpora along with the evaluation measures in this field.

Font et al. [19] performed a comparative analysis of nine different countermeasures with cross-corpus evaluation of the countermeasures. However, the authors only discussed the performance of the countermeasures against replay attacks due to the nature of the dataset used for the study. Although Sahidullah's (2016) [70] study and Font's (2017) [19] paper each provided a comprehensive analysis of the replay attack, other types of spoofing attacks that could severely disrupt an ASV system were ignored. Sahidullah et al. surveyed four

different types of spoofing attacks, spoofing procedures, and countermeasures the next year, in 2019 [69], summarized some spoofing challenges and presented countermeasures. However, a thorough examination of all aspects of the ASV system was still wanting. For instance, [70] and [19] focused only on front-end countermeasure design, and concentrated on the SS, VC, replay, and mimicking attack types, alone where speech corpora, protocols, classifier, and evaluation metrics all contributed equally to the construction of an ASV system. These were absent from the [69] study.

Aside from [79], most of the articles did not present a detailed taxonomy of modern voice presentation attack detection (PAD) methods, where the work was structured based on identified elements. This survey, along with [56] broadened the classification of relevant work and built on the taxonomy from the most recent work on voice PAD. The significant contribution of these articles was the presentation of trends and analyses of voice PAD that were absent in the other survey articles. Mittal et al. (2021) described the computation mechanisms of traditional and modern speech feature extraction, as well as datasets and a combination of various front- and back-end techniques. None of these articles, however, performed a fair experimental analysis of the existing state-of-the-art countermeasures. In addition, existing reviews and survey articles also lack cross-corpus evaluation in order to show the generalizability of existing solutions. The comparison of the existing reviews and surveys is presented in Table I.

Existing spoofing countermeasures have employed diverse features and classifiers and, typically, have had their performance evaluated on only one of several datasets, i.e., ASVspoof 2017, ASVspoof 2019, and ASVspoof 2021, using different metrics. With no standardization, it is difficult to declare a single countermeasure as the unified method that works best. Thus, there exists a need to provide a thorough analysis of existing countermeasures in order to show which method is the best fit for a certain scenario, i.e., spoofing type, dataset, single or multi-hop attack, etc. To the best of our knowledge, no comparative analysis work on multiple voice spoofing attacks, including single and multi-order attacks, has ever been presented. Moreover, existing studies have ignored the important aspect of cross-corpora evaluation, which is crucial to the evaluation of the generalized nature of the countermeasure. To address these challenges for the existing countermeasures, we present a comparative analysis of various spoofing countermeasures. Motivated by discussions, the survey proposed in this paper investigates significant contributions to the ASV system development chain. More specifically, the main contributions of our work are:

- We present a baseline survey of voice spoofing countermeasures, involving diverse factors, for the benefit of upcoming anti-spoofing systems.
- We present a comprehensive analysis of existing state-of-the-art voice spoofing attacks, countermeasures (hand crafted features, along with deep learning and end to end solutions), publicly available datasets, and the performance evaluation parameters used in voice spoofing.
- We present an experimental analysis of state-of-the-

art different countermeasures, on several classifiers, for spoofing detection performance using the VSDC, ASVspoof 2019 and 2021 datasets, and evaluate the results across corpora.

- We cover one of the key limitations of cross corpus evaluation and generalization mentioned by the in-field researchers and evaluate the performance of the featured countermeasures on three large-scale publicly available and diverse datasets, i.e., ASVspoof 2021, 2019 and VSDC, in terms of min-tDCF and EER. The countermeasures are tested against four different machine learning and deep learning classifiers.

## II. Literature collection and selection criteria

In this survey, we review existing research papers that focus on techniques for detecting and countering audio spoofing attacks. A detailed description of the approach and protocols employed for the review is given in Table II. We also display trends in spoofing attacks and countermeasures year over year by examining Google Scholar papers released in the previous 6 years (2015–2022). Figure 1 depicts the course of the published study. The rest of the paper is structured as follows: Section 3 provides current voice spoofing countermeasures. Section 4 discusses voice spoofing attacks as well as adversarial attacks, and Sections 5 and 6 discuss publicly available data sources and performance evaluation strategies. Section 7 describes the acoustic features and classifiers used in a comparative analysis of the countermeasures. Sections 8 and 9 present the experimental setup and analysis, as well as a discussion of limitations and future trends. Finally, Section 10 provides our conclusions and motivation for future work.

## III. Voice Spoofing Attacks

ASV systems, are vulnerable to a variety of voice spoofing attacks, i.e., physical-access (PA) and logical-access (LA) attacks. These and other possible attacks on ASV systems are graphically illustrated in Figure 2. PA attacks are easier to conduct, compared to LA attacks, due to the availability of modern high-quality recording devices. A replay attack can be generated using a device as simple and commonly available as a mobile phone, compared to the robust and sophisticated AI voice cloning algorithms required for LA attacks. Hence, voice replay attacks are more commonly employed to spoof ASV

systems in order to accomplishing an attacker's objectives e.g., to gain entry to someone's home. Moreover, we have demonstrated in our earlier work [6] that voice replays can be generated not only in a single order but also in multi-order attacks, by using different smart speakers such as Amazon Alexa and Google Home.

Voice spoofing countermeasures, also known as presentation attack detection (PAD) systems, aim to detect the following attacks (see also Figure 3):

- Direct attacks can further be broken down into device artifacts and algorithmic artifacts, which are **A)** Physical access (PA) and **B)** Logical Access (LA) scenarios, respectively.
  - PA attacks occur when the samples are applied as an input to the ASV system through the sensor, such as a microphone.
    * Replay attacks are one example of ASV PAs. Replay attacks occur when the audio of an authorized user is recorded and played back in order to deceive an authentication system. Replay attacks are the easiest to attempt, as they do not require any special technical knowledge, and can easily deceive ASV systems. Because replayed audio includes background and/or other noise, however, it can be detected.
    * An impersonation attack occurs at the microphone level, where a person changes how they talk to mimic the speech characteristics of a legitimate



**Fig. 1:** Number of papers in the area of voice Spoofing Attacks and Countermeasures by yearly publication.

**TABLE I** Comparison of the existing survey and review papers

| Paper | Presentation Attacks | | | Countermeasures | | | Integrated ASV | Experimental Analysis | Cross Corpus |
|---|---|---|---|---|---|---|---|---|---|
| | Replay | VC/Synthesis | Adversarial | Hand Crafted | Deep Learned | End-to-End | | | |
| [99] | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| [70] | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ |
| [19] | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ |
| [64] | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| [69] | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| [79] | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| [56] | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| [36] | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| [Ours] | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

**TABLE II** Literature collection protocol.

| Preparation Protocol | Description |
|---|---|
| Purpose | • To identify current state-of-art voice spoofing attacks and countermeasures. <br> • To critically compare both popular approaches and the features used in voice spoofing countermeasures. <br> • To investigate open challenges that still need investigation in the domain of voice spoofing countermeasures, and paths forward. |
| Sources | Google Scholar, IEEE explore , Springer Link |
| Query | Queries were used on the data sources above for collection of sources: Audio Spoof Countermeasure/ Spoofed Audio Detection/ Audio Synthesis Detection/ Audio Replay Detection/ Audio Spoofing Countermeasures/ Automatic Speaker Verification/ Secure ASV System/ ASV spoofing/ Presentation Attacks/ Voice Conversion/ Adversarial Attacks on ASV/ Voice Replay Attacks/ Anti Spoofing Countermeasures/ Voice Spoofing Detection. |
| Method | Literature was categorized as follows: <br> • Audio spoofing detection and countermeasures based on traditional and deep-learned methods, as well as handcrafted features. <br> • Existing voice spoofing attacks, along with adversarial attacks and their taxonomy. <br> • An extensive examination, testing and comparison of existing voice spoofing countermeasures using single and cross corpus evaluation. <br> • Detailed discussion of the research gap, issues, limitation, and future trends in voice spoofing detection and countermeasures. |
| Size | A total of 150 papers were retrieved using this query mentioned above from the method until the search was discontinued on 04-20-2022. We selected literature that was relevant to the subject of audio spoof detection and excluded those that were not relevant to this subject or were white papers/articles. |
| Inclusions and Exclusions | Preference was given to peer-reviewed journal papers and conference proceedings articles. In addition, articles from the archive literature were also taken into account. |

user. An attacker can fool an ASV system if the impostor's natural voice has similar features.

– LA attacks occur when the audio samples bypass the sensor and are injected directly into the model. Spoofing attacks based on logical access include **C)** voice conversion, **D)** speech synthesis, and E) Adversarial attacks.

* Voice conversion uses an imposter's natural voice to generate artificial speech in order to match the targeted speaker's voice. Attacks are usually created by other models to fool a specific ASV system. Machine learning has allowed the mapping of speech features between speakers to be accurate, and is now computationally efficient enough to make ASV systems vulnerable to these types of attacks. However, these attacks can still be detected because they are not a perfect match to genuine audio.

* Speech synthesis attacks, also known as deepfake attacks, are similar to voice conversion but use text as an input to the model in order to generate a voice clip similar to a targeted speaker in an effort to fool an ASV system. Models that generate accurate speech features can be trained using a small data set of recorded audio.

• Indirect attacks include adversarial attacks, where the audio signal remains unchanged while the attacker modifies the signal's properties during ASV processing. The recent work on adversarial attacks is discussed later in the paper.

*A. Proposed countermeasures to Adversarial Attacks*

Adversarial attacks, such as the addition of Gaussian noise, can create small perturbations in audio samples and pose a threat to ASV anti-spoofing models, affecting accuracy and safety. The method described in [96] uses self-supervised learning to leverage the knowledge of unlabeled data, which significantly improves performance. They create high-level representations, extracted by the self-supervised model, and a layer-wise noise to signal ratio (LNSR) in order to quantify and measure the effectiveness of deep models in countering adversarial noise. During training, different masking methods are applied to 15% of the audio samples to create a model that reconstructs audio. This defense is passive, i.e., not a proactive defense. Their proposed method, called Mockingjay, filters the audio. This model was only tested on the Projected Gradient Descent (PGD) and Fast Gradient Sign Method (FGSM) attacks but not on other attacks. Also, as epsilon (the strength of an attack) increases, the defense starts to weaken and eventually fails. Both white-box and black-box scenarios were tested.

Another defensive method, [95], uses spatial smoothing, filtering, and adversarial training. This method has only been tested against PGD attacks and may perform poorly under other, perhaps stronger, attacks because different attacks leave different patterns on the sample. These defense methods improved the robustness of spoofing countermeasures against PGD attacks, but they are still limited in their overall defensive capability because adversarial training becomes weaker against unfamiliar attacks. The work of [98] creates a model that fits the distribution of genuine speech, where it takes genuine speech as the input and generates "genuine" speech, or output with the same distribution as the genuine speech, as the output. However, for spoofed speech, it will generate very different output that amplifies the difference in the distribution compared to genuine speech. Therefore, they propose a genuinization transformer that uses genuine
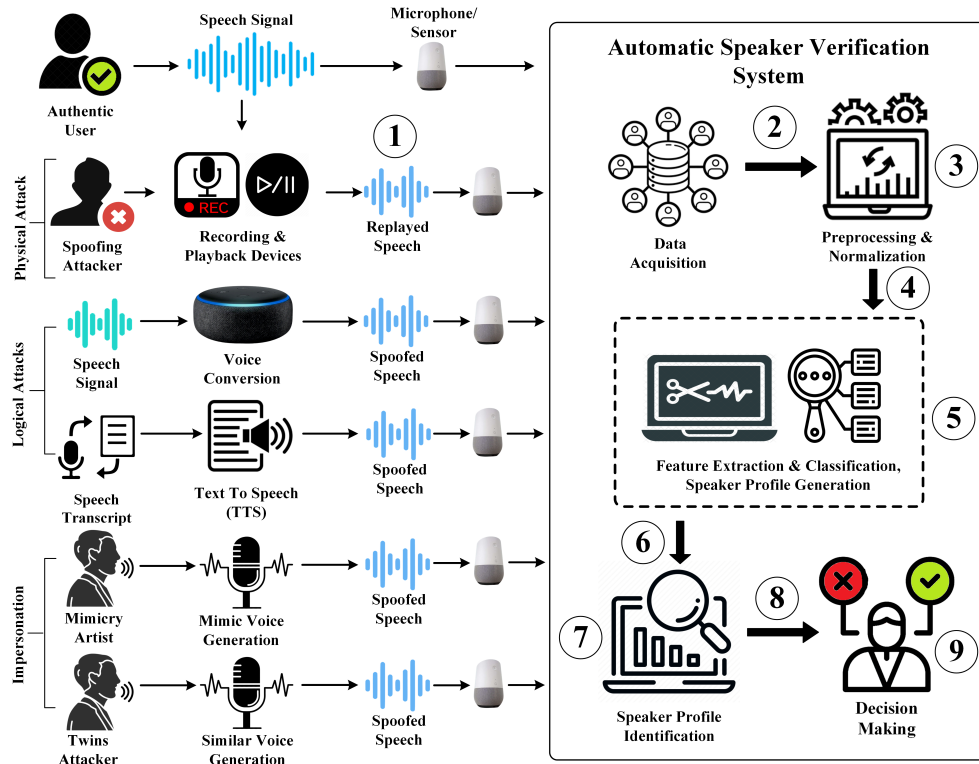
**Fig. 2:** Existing Threats to Automatic Speaker Verification Systems.

speech features with a convolutional neural network (CNN). The genuinization transformer is then used together with a Light Convolutional Neural Network (LCNN) system for the detection of synthetic speech attacks. It was trained on the ASVspoof2019 dataset with Constant Q Cepstral Coefficients (CQCC) and Linear Frequency Cepstral Coefficients (LFCC) as features, and achieved an EER of 4.07% and a min-TDCF of 0.102. However, replay attacks were not done, and this model was not tested across corpora. This paper, [105], introduces an inversion module that derives four features: inverted constant-Q coefficients (ICQC), inverted constant-Q cepstral coefficients (ICQCC), constant-Q block coefficients (ICBC), inverted constant-Q linear block coefficient (ICLBC). The first two use conventional Discrete Cosine Transform (DCT) and the latter two use overlapped block transform. These features for synthetic speech detection are evaluated using the ASVspoof 2015, noisy ASVspoof 2015, and ASVspoof 2019 logical access databases, and achieve an EER of 7.77% and a min-TDCF of 0.187. Suthokumar et al, in [74], compare how genuine and spoofed speech varies across different phonemes, and identifies certain phonemes (fricatives, nasals, stops, and pauses) which are more informative in the detection of replay attacks. The paper created four different fusion scoring methods to incorporate phonetic information using phoneme-specific models, and achieved an EER of 6.18% on the ASVspoof2017 V2 dataset. However, this model has to be fused with a phenome independent model for better results. For logical scenarios, this paper uses feature extraction, a densely connected network, a squeeze and excitation block, and a feature fusion strategy. For physical scenarios, their

method consists of feature extraction, a multi-scale residual network, an SE block and weighted average strategy. Both methods for logical and physical scenarios combine physical and physiological characteristics to resist spoofing attacks.

Adversarial attacks may also pose a threat to ASV systems. According to [27], ASV systems are prone to adversarial attacks, which may reduce the accuracy of such systems by up to 94%. Attacks can range from simple Gaussian noise to more advanced attacks created in a targeted white-box setting. These perturbations, imperceptible to humans, may cause audio classification and ASV systems to fail completely.

Adversarial training is another way to reduce the effectiveness of an adversarial attack, as done by [27]. In this paper, adversarial training– when adversarial samples are included in the training dataset– and Adversarial Lipschitz Regularization (ALR) are used to reduce the impact of adversarial attacks. ALR is based on a function that ignores small changes in the input. Adversarial training is beneficial for defense when the attacks are similar to each other, but not all attacks will be similar. This method starts to fail when introduced to different types of attacks, some stronger than others, and this failure is evident when they are evaluated. The work of [92] uses neural Vocoders to re-synthesize audio and finds differences between the ASV scores of an original audio sample and a re-synthesized audio sample. This method detects adversarial samples by first purifying the sample, i.e., removing adversarial noise while generating the genuine waveform with reduced distortion. This method is beneficial because it does not need to know the attack algorithm used. This model is tested purely on a Basic Iterative Method (BIM) adversarial attack,
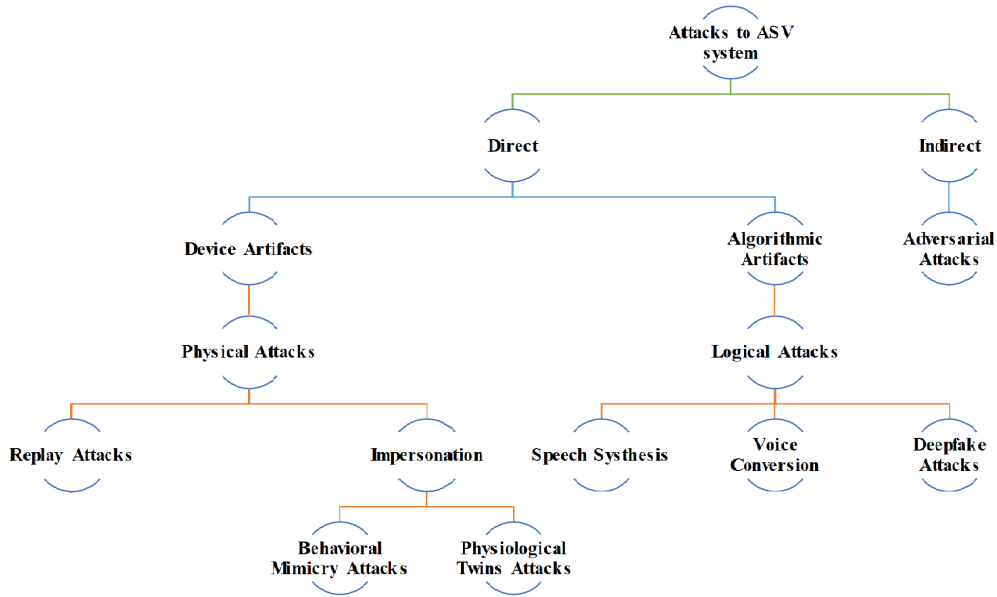
**Fig. 3:** Taxonomy of the Voice Spoofing Attacks.

which extends an FGSM attack where FGSM is performed multiple times with a small step size. Adversarial attack types exhibit different effects, so judging performance on many types of attacks is difficult. Further, only one type of neural vocoder was used. Vocoders, when used to reconstruct phase information, inevitably introduce noise and distortion. The authors claim that this noise is actually beneficial, but this may not be entirely accurate. In addition, self-supervised learning-based models are more common when it comes to detecting and purifying adversarial samples. Self-supervised learning may be used to reconstruct audio, which is beneficial when it comes to adversarial attacks on audio.

The work of [93] uses self-supervised learning-based models for adversarial defense on ASV systems by using Trans-former Encoder Representations from Alteration (TERA) models. The authors use self-supervised learning in a similar way to Mockingjay, which was a considerable influence on this paper. This method is beneficial because it does not require knowledge about how the adversarial samples were generated, as it purifies the samples. The authors use Gaussian, mean, and median filters for audio samples. Similarly, the work of [94] uses Self Supervised Learning Models (SSLMs) to remove superficial noise from the inputs, and reconstructs clean samples from the interrupted ones. They have two modules; one module purifies audio, while the other compares the original and purified samples in order to detect adversarial samples.

**TABLE III** A Comparison of Adversarial Defensive Methods for Anti-Spoofing Models

| Paper | Technique | Dataset | Adversarial Testing | Key Performance | Limitations |
|---|---|---|---|---|---|
| [96] | Mockingjay with SSL | ASVspoof 2019 LA and LibriSpeech | FGSM and PGD | 90% accurate | limited to 2 attacks testing |
| [95] | Spatial Smoothing and Adversarial Training | ASVspoof 2019 LA | PGD | 92.40% (SENet), 98.60% (VGG) | Limited to 1 attack only |

**TABLE IV** Comparison of Adversarial Attack Defense for ASV Systems

| Author | Technique | Dataset | Tested Adversarial Attacks | Best Evaluation Performance | Limitations |
|---|---|---|---|---|---|
| [27] | Adversarial Lipschitz Regularization | Librispeech | PGD, FGSM, Carlini, Wagner | 73% adversarial accuracy | Fails with stronger attacks |
| [92] | Neural Vocoders | Voxceleb1 and Voxceleb2 | BIM | AUC: 99.94 | Single Attack Testing |
| [93] | TERA models | Voxceleb1 | BIM | 22.94% EER | Only tested on 1 type of attack |
| [94] | SSLMs | Voxceleb1 | BIM, FGSM, JSMA | R-Vector 14.59% GenEER X-Vexctor 11.29% GeEER | Degrades performance |

## IV. Taxonomy of Voice Spoofing Countermeasures

This section provides a detailed analysis of existing state-of-the-art voice spoofing techniques adopted to detect spoofing attacks as well as countermeasures developed to combat those attacks. The countermeasure taxonomy illustrated in Figure 4 includes a comprehensive classification of the countermeasures.

### A. Conventional Handcrafted Spoofing Countermeasures

The countermeasures for audio spoof detection are classified into two categories: conventional handcrafted features and enhanced deep learning solutions. For voice spoofing detection, handcrafted features are commonly employed, i.e., MFCC [112], CQCC [83], and others, with conventional machine learning classifiers, e.g., Gaussian Mixture Models (GMMs) [7], Support Vector Machine (SVM) [87], etc.

In [65], the author's novel acoustic features were derived from the frequency-warping and block transformation of filter bank log energies. A two-level Mel and speech-based wrapping was used with an Overlapped Block Transformation (OBT) to extract the conventional and inverted classification features. Although the proposed features achieved a 0.99 classification accuracy on the ASVspoof 2015 corpus development section, they were only applicable when the spoofing attacks were known in advance. Aside from this, the performance of the proposed approach was not evaluated and reported in the presence of unknown adversarial attacks. Further, the approach was not validated across corpora due to a lack of open-source datasets. In the ASVspoof2015 challenge, the unified articles concentrate on countermeasures for speech synthesis and voice conversion spoofing attacks. Conversely, replay attacks pose the greatest threat to ASV systems. These attacks involve the playback of recordings, acquired from registered speakers, in order to generate false authentication tokens, and may be efficiently carried out using common consumer devices.

ASVspoof2017, the second in the series of challenges, focused on the development of replay attack countermeasures. In [39], Kinnunen et al. mainly followed this line of thinking, utilizing front-end features, i.e., CQCC, MFCC, PLP, RFCC, Inverted Mel Frequency Cepstrum Coefficients (IMFCC), LPCC, SCMC, SSFC, and HPCC, with a baseline GMM-UBM classifier. The performance of these countermeasures was evaluated on the ASVspoof2017 dataset, where the proposed systems achieved an EER of 31.5% for genuine vs. replay and 1.8% EER for genuine vs. zero-effort imposters. In addition, the average detection EER of all main submissions was 25.91%, whereas the best single system result had an average detection EER of 6.73%. When these results were compared to those from the prior challenge (ASVspoof2015), the detection of replay attacks was found to be more complex than detecting speech synthesis and voice conversion spoofing attacks. In [70], Sahidullah et al. present the first comparative evaluation of six countermeasures (CMs) and their integration with automatic speaker verification using the ASVspoof 2015 dataset. The CMs contain MFCCs, IMFCCs, LFCCs, Cochlear Filter Cepstral Coefficients (CFCCs), CQCCs, and Gammatone Frequency Cepstral Coefficients (GFCCs) as front-end cepstral features, which are then coupled with GMM-UBM and i-Vector classifiers for back-end speaker classification. According to the results presented in this paper, the CM using CQCC features and the fusion of all six CMs has the greatest potential to detect spoofing attacks in the context of spoofing attacks. More significantly, the fusion of CM and CQCCs achieves the lowest EER of 0.02%.

Cascading integration of ASV and CMs significantly decreases the False Acceptance Rate (FAR). Although the significance of CQCC features is demonstrated in this study, the presented system is not evaluated against other forms of spoofing attacks, e.g., replay, deepfake, adversarial, or others. Further, countermeasures are not tested cross-corpus. The results of [70] demonstrate that an ensemble of acoustic features, together with machine learning models, can produce more accurate results than individual classifiers. Consequently, in [30], the author presents an ensemble of the acoustic features of CQCCs, along with other classical features, i.e., MFCCs and Perceptual Linear Predictive (PLP) features. The author proposes an ensemble classifier set that contains numerous GMMs, i.e., Gradient Super Vector-Boosting Decision Tree (GSV-GBDT) and GSV-Random Forest (GSV-RF) classifiers. The experimental results show that the presented ensemble system significantly outperforms the baseline GSV-SVM system. Specifically, using a baseline GSV-SVM classifier, this approach achieves an EER of 10.4% on CQCC features, 27.4% on MFCC features, and 37.0% on PLP features. In contrast, a minimal EER of 9.5% is achieved by employing the ensemble model. Similarly to [65], the biggest limitation of [30] is that it has only been tested against a single replay attack. Other types of spoofing attacks, i.e., speech synthesis, voice conversion, and the most current deepfake and adversarial attacks, are not reviewed, and there is no cross-corpus examination of the techniques. In [4], the authors examine the robust audio features, comprised of handcrafted and auto-encoder-based learned features, to identify replay spoofing attacks. The handcrafted features employed in this study are CQCCs, LPCCs, IMFCCs, Rectangular Filter Cepstral Coefficients (RFCCs), Sub-band Centroid Frequency Coefficients (SCFCs), and Sub-band Centroid Magnitude Coefficients (SCMC), as well as spectrogram features. The paper uses an auto-encoder to learn a dense representation of all of the features. A conventional GMM, along with a Universal Background Model (UBM), is used as the baseline system to examine the performance of the handcrafted and encoder-based features. The integrated fused models, based on existing audio and machine-learned features, achieve comparable results, with an identical error rate (EER) of 12.0%. In particular, the handcrafted CQCC features outperform all other features, with an EER of 17.5%. The best encoder-based feature observed was a spectrogram with a minimum EER of 20.2%. Although the coupling of handcrafted features with an auto-encoder-based system surpasses the state-of-the-art, the presented system is only tested against replay attacks. In the presence of other types of spoofing attacks, e.g., voice conversion and adversarial attacks, the system's performance tends to vary.
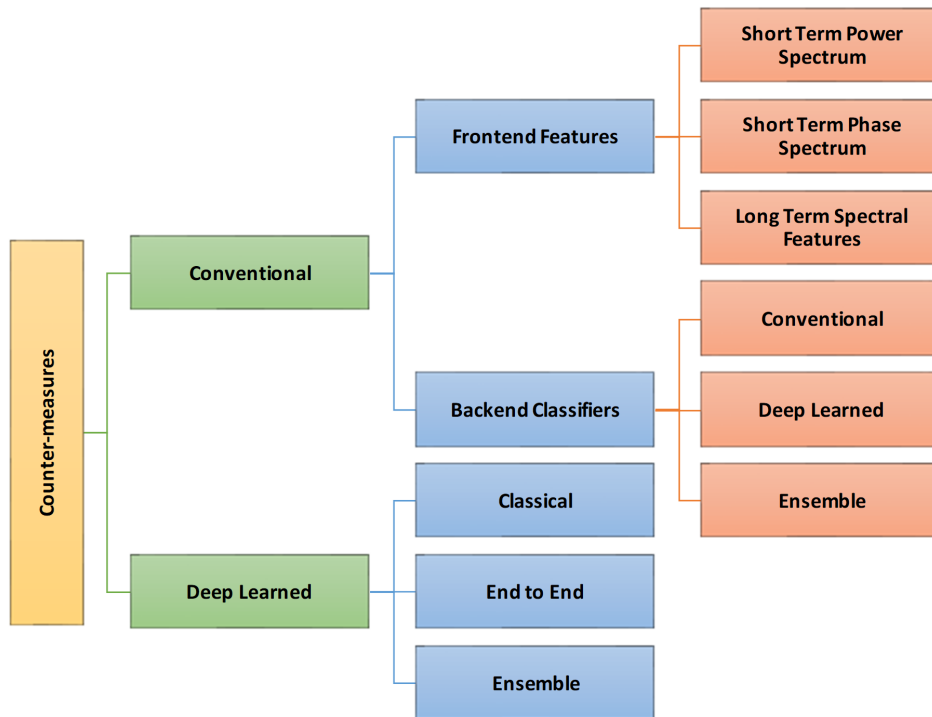
**Fig. 4:** A detailed taxonomy of existing voice spoofing countermeasures.

Rather than using standard standalone short-term power spectrum coefficients, i.e., MFCCs, LFCCs, CQCCs, etc., to identify the spoofing attacks, the authors of [81] introduced phase-based Teager Energy Operator (TEO) features. It was discovered that the TEO phase features gave information that is complementary to the information provided by the more commonly used CQCC, MFCC, and LFCC feature sets. Although the TEO phase features were unable to perform well alone, fusion with traditional features enhanced the accuracy of the spoofing detection system. The results demonstrated that the standalone spoof detection systems, developed with the TEO phase, MFCC, and LFCC, achieved an EER of 31.34%, 34.02%, and 16.80%, respectively, whereas, when the TEO phase feature set was fused with the CQCC, MFCC, and LFCC feature sets, the EER was lowered by 0.18%, 2.74%, and 1.41%, respectively. This improvement in system performance showed the influence of TEO phase information on the spoof detection system. However, the provided solution was only stated to be resistant to replay spoofing attacks, and had no cross-corpus validation.

In accordance with [81], the authors of [63] present phase spectrum and multi-resolution wavelet features, in addition to the commonly used front-end MFCC features. This study combines MFCCs with Mel-Frequency Principal Coefficients (MFPCs), CosPhase Principal Coefficients (CosPhasePCs), and Mel Wavelet Packet Coefficients (MWPCs) to provide a reliable and robust defense against spoofing attacks. The experimental results indicate that applying principal component analysis (PCA) to MFCCs results in a considerable EER improvement over MFPC features for all spoofing methods. However, the MFCC features prove inferior in comparison with other front end features. In contrast, the implementation

of CosPhasePC modestly decreases EER in comparison to the MFPC features. Although the multi-resolution wavelet transform features outperform the state-of-the-art, achieving 0.05% EER for all known attacks, the proposed framework has not been tested and reported cross-corpus and against unknown attacks. Accordingly, following the significant performance of phase-oriented features, the author sheds light on the importance of acoustic front-end features and introduces a novel detection mechanism by modeling replayed speech as a convolution of original speech [80]. Also in [80], the author proposes a novel feature set, Magnitude-based Spectral Root Cepstral Coefficients (MSRCC) and Phase-based Spectral Root Cepstral Coefficients (PSRCC), which outperforms the baseline system (CQCC) and provides a 29.18% EER on the evaluation set of the ASVspoof 2017 challenge database. With the GMM back-end classifier, the front-end features MSRCC and PSRCC respectively produce 18.61% and 24.35% EER. Conversely, with convolutional neural network (CNN) back-end classifiers, MSRCC and PSRCC obtain 24.50% and 26.81% EER, respectively. In addition, the score-level fusion of MSRCC and PSRCC results in 10.65% and 17.76% EER using GMM and CNN classifiers, respectively. These findings suggest that the proposed feature sets of MSRCC and PSRCC capture complementary information. However, the cross-corpus examination was not reported.

The research in [8] shows the effectiveness of augmenting genuine training data in the simulation of replay spoofing attacks. The author presents replay spoofing countermeasure systems that improve the CQCC-GMM baseline with score level fusion. Instead of using CQCC, the author used spectrograms as inputs to analyze end-to-end feature representations. Finally, the author replaced the subsequent GMM classifier

with a Fully-connected Deep Neural Network (FDNN) and a Bidirectional Long-Short-Term Memory neural network (Bi-LSTM). The results of the experiments show that this data augmentation technique can significantly increase the system's performance. In particular, the baseline CQCC-GMM model achieved an EER of 22.29%, while the DA-CQCC-GMM model obtained an EER of 19.18% and the fused system achieved an EER of 16.39%. Although the end-to-end FDNN and Bi-LSTM-based systems perform well for replay detection, other sorts of current deepfake and adversarial attacks may have led the system to perform inadequately. Furthermore, the provided system was not evaluated in the presence of unknown attacks and audio recordings across corpora.

Yang et al. [106] developed a Low-Frequency Frame-wise Normalization (LFFN) technique to capture replay distortions. LFFN was combined with Constant-Q Transforms (CQT) to extract two features: Constant-Q Normalization Segmentation Coefficients (CQNSC) and Constant-Q Normalization Cepstral Coefficients (CQNCC). This approach performed well on the ASVspoof 2017 version 2.0 dataset, with an EER of 10.63% on CQNSC and 10.31% on CQNCC features. The novel Acoustic Ternary Patterns-Gammatone Cepstral Coefficient (ATP-GTCC) feature was proposed in [53] to help develop a lightweight model for single and multi-order replay attack detection. An SVM classifier was used to better capture the harmonic distortions found in multi-order replay samples, while ATP-GTCC was employed as a front end feature. This model achieved an EER of 0.6% and 1% on the Voice Spoofing Detection Corpus (VSDC) and the ASVspoof 2019 dataset, respectively. This system exceeded earlier state-of-the-art techniques in terms of replay detection performance and efficiency. However, the given approach is limited to single-order and multi-order replay attacks, and has not been tested against other types of attacks, including voice conversion, deepfake, and adversarial attacks.

In [100], Wu et al. use a common benchmarking database to analyze the vulnerability of text-dependent speaker verification systems to the replay attack and propose an anti-spoofing technique to protect the ASV systems. The main idea behind this spoofing detection approach is to use a similarity score to determine whether the provided sample is identical to any previously stored voice samples. A back-end HMM classifier is trained using the front-end MFCC, LPCC, and spectro-temporal features. Experiments on the RSR2015 database reveal that in the presence of a replay attack, the EER and FAR, which were identical at 2.92%, increased to 25.56% and 78.36%, respectively, proving the speaker verification system's vulnerability to replay attacks. In the presence of replay spoofing attacks, the proposed spoofing countermeasure lowers the FARs from 78.36% and 73.14% to 0.06% and 0.01% for male and female samples, respectively. In [24], the authors use a feature fusion of GTCC, MFCC, Spectral Flux, and Spectral Centroid for input audio presentation. This countermeasure successfully detects multiple types of logical access attacks and classifies the cloning algorithms used to produce the synthetic speeches. It achieves an EER of 3.05%, compared to baseline models that achieve an EER ranging between 5.06% and 9.57%. The authors of [3] present a

voice spoofing countermeasure using ELTP-LFCC features and a Deep Bidirectional LSTM to combat TTS synthesis and converted voice samples in LA attacks. In this paper, ELTP is further fused with LFCC to better capture the characteristics of the vocal tract speech dynamics of both bonafide voice and cloning algorithm artifacts. On the diverse ASVspoof 2019-LA dataset, performance evaluation yields an EER of 0.74% and a min-TDCF of 0.008%. However, the presented system is only tested on one dataset.

Rigorous experimentation was performed to illustrate the significance of the proposed countermeasure for detecting LA-based voice spoofing attacks. In the literature, frequency or amplitude modulation-based features were also explored for voice spoofing detection. Gunendradasan et al. [23] used Spectral Centroid Deviation (SCD) features to develop a replay attack detection system. The Spectral Centroid Frequency (SCF) and Spectral Centroid Magnitude Coefficient (SCMC) features were extracted from the same front-end as SCD and used as complementary features. These features were used to train a GMM classifier for replay attack detection. This method was evaluated on the ASVspoof2017 dataset and provided a 60% improvement in EER as compared to the ASVspoof CQCC baseline model.

Existing methods have also used various handcrafted features to address both physical (replay) and logical access (clone, voice conversion) spoofing. Monteiro et al. [57] proposed a system that can detect both logical and physical access audio attacks against ASVs via an ensemble of three different models. The baseline systems of LFCC-GMM and CQCC-GMM obtained an EER of 8.09% and 9.57%, respectively, on the evaluation set. Although this system achieved an EER of 9.87% for voice cloning attacks, for replay attacks, the baseline systems achieved an EER of 11.04% and 13.54%, whereas the proposed system achieved an EER of 1.74%. Zhizheng et al. [104] proposed an ASV spoofing and countermeasures initiative with more focus on speech synthesis and voice conversion spoofing attacks (presentation attacks). The authors also described the post-evaluation results of the ASV systems achieved by the fusion of several spectral features, i.e., MFCC, CFCCIF, and MFPC, as well as the back-end classifiers, such as GMM, SVM, SVM-Fusion, and SVM with i-Vectors.

Some methods also exploit high-frequency components in order to capture the traits of bonafide and spoofed signals. The work in [90] captured the high-frequency content by using the inverted-MFCC (IMFCC), LPCC, LPCCres, CQCC, MFCC, and Cepstrum features. This hybrid feature representation was then used to train a GMM for the classification of bonafide and spoofed samples. Due to the increased feature computation cost, this method was not suitable for local deployment on resource-constrained voice-controlled systems. This paper, [22], also highlighted the idea of analyzing the high-frequency band for replay spoof detection. For this purpose, transmission line cochlea amplitude modulation and transmission line cochlea frequency modulation features were employed to train a GMM for replay attack detection. This method was evaluated on the ASVspoof2017 dataset and achieved an EER of 7.59%. Although this method [22] performed only slightly better than the baseline model, amplitude modulation-based features have

a high computational cost because they take more than twice the amplitude frequency to modulate the signal.

In the literature, non-voice segments have also been explored to capture the distortions of playback speech. Saranya et al., in [72], analyzed the channel and reverberation information from non-voiced segments of the input audio for replay spoofing detection. A voice activity detector was used to determine the non-voiced segments, then a hybrid feature vector comprised of CQCC, MFCC, and Mel-Filter bank-Slope was employed to capture the remnant vocal tract information in the non-voiced segments of the input audio. These features were later used to train a GMM for the classification of bonafide and replayed audio samples. The performance of this method was evaluated on the ASVspoof2017 dataset, and showed an improvement of 37% in EER over the baseline method. In the work of [105], the authors introduced an inversion module to derive four features: ICQC, ICQCC, ICBC, and ICLBC. Two of them used conventional DCT, and the other two used overlapped block transforms. These features for synthetic speech detection were evaluated using the ASVspoof 2015, noisy ASVspoof 2015, and ASVspoof 2019 logical access datasets, and achieved an EER of 7.77% and a min-tDCF of 0.187. Despite the fact that they were tested on multiple datasets, due to aadvances in spoofing attacks, i.e., adversarial or unknown attacks, they may not perform well in a real-world scenario.

## B. Discussion

We performed a survey to determine the most frequently used features for spoofing audio and speech countermeasures. Each handcrafted feature was carefully designed and has its own dependencies to differentiate between genuine and spoofed voice samples. Despite the fact that each spoofing attack adds a distinct distortion to the speech sample in order to mislead an ASV system, some recently published crafted features worked very well in the classification of physical and logical attacks. Although the reported hand-crafted features have been evaluated on several datasets, including ASVspoof2015, ASVspoof2017, ASVspoof2019, and others, they are dependent on the data source used to train and classify the spoofed and bonafide speech samples. There is a significant need for a standardized, unified dataset to train and validate the accuracy of the published feature sets in order to choose the best among them. Furthermore, the existing features have only been tested on training datasets so far, with no cross-corpus assessment. Consequently, the list of features we have chosen covers the bulk of the power, frequency, duration, amplitude, and phase spectrum in order to distinguish between bonafide and spoofed voice samples. Moreover, we offer a cross-corpus examination of the features and test performance on standard datasets. The experiment and results section provides the details of the performance analysis on the selected speech samples.

## C. Deep learning based Countermeasures

Recent years have witnessed a rise in the use of deep learning approaches to prevent audio spoofing attacks. How-

ever, despite the benefits of handcrafted features, these approaches were computationally complex and required powerful computing resources. Inspired by the great success of deep learning in automatic speech recognition, deep neural network (DNN) based systems were developed for spoofing detection for the first time in 2015. In [10], a novel, simple model for detecting spoofing attacks on a speaker verification system was developed. The presented model was used to extract key features from audio samples and construct compact, abstract, and resilient deep data representations. A spoofing-discriminant network was used in the training of spoofing algorithms. The proposed network then computed the s-vector, which is the utterance level average of the final hidden layers. Finally, Mahalanobis distance, along with normalization, was used, in conjunction with computed s-vectors, to detect spoofing attacks. This model achieved the 3rd position in the first spoofing detection challenge, ASVspoof 2015. In particular, the proposed model attained an overall minimum EER of 2.281%. Specifically, the presented system obtained an EER of 0.046% for known and 4.516% for unknown replay attacks. Although the proposed system performed well in the ASVspoof 2015 competition, it has not been examined across corpora. Furthermore, when spoofing techniques other than the replay attack are introduced into the system, the performance of the proposed technique is taken into account.

Existing methods also employed handcrafted features to train deep neural networks for spoofing detection. In [60], high-frequency regions were analyzed by proposing the High Frequency Cepstral Coefficients (HFCC) feature and using a back-end DNN model to classify bonafide and spoofed (replay) audio. When comparing HFCC features to CQCC features using a baseline GMM back-end, HFCC outperformed CQCC in the development and evaluation sets of ASVspoof2017. HFCC received an EER of 5.9% in the development set and 23.9% in the evaluation set, while CQCC received 11.0% EER in the development set and 24.7% in the evaluation set. An attentive filtering network system was proposed in [45] to detect replay attacks against ASV systems, and the proposed method enhanced feature representations in both the frequency and time domains. This system obtained an EER of 6.09% and 8.54% on the development and evaluation sets of ASVspoof 2017, respectively, on a system comprised of two Attentive Filtering models (one using sigmoid and the other using softmax as their activation functions). This method achieved better results than the CQCC-GMM baseline model, which obtained an EER of 12.08% and 29.35%. However, the combination of feature extraction and the model's architecture required significant computational resources.

In [25], segment-based linear filter bank features, along with an attention-enhanced DenseNet-BiLSTM model, were proposed for replay attack detection. These features were extracted from the silent segments of the audio signal. Triangle filters were used to examine noise in the high frequency bands $3 - 8kHz$. The baseline system of CQCC-GMM obtained an EER of 2.36% on the development set and an EER of 8.42% on the evaluation set of ASVspoof2017, while features of this method, when used with a GMM, obtained an EER of 1.8% and 7.92% on the development and evaluation sets,

respectively. One limitation of this method was the sensitivity of the signal-to-noise ratio during segmentation of the silence and voiced components, where low SNR resulted in false segmentation.

Even though the work of [98] uses LPS features, the enhanced modification was before passing them on to an LCNN. The authors create a model that fits the distribution of genuine speech, i.e., one that takes genuine speech as the input and generates genuine speech as the output. However, if the speech is spoofed, the output will be very different. They propose a genuinization transformer that uses genuine speech features with a convolutional neural network (CNN). The genuinization transformer is then used with an LCNN system for the detection of synthetic speech attacks. The model achieves an EER of 4.07% and a min-tDCF of 0.102 on the ASVspoof2019 dataset, using CQCC and LFCC features. However, this method has not been tested across corpora. The work of [52] proposes a Conv1D Resblock with a residual connection, which allows the model to learn a better feature representation from raw waveforms. They find that feeding a raw waveform directly into a neural network is adequate. Only the ASVspoof2019 dataset is used, and it achieves an EER of 2.98% and a min-tDCF of 0.082. The work of [76] is an extension of a previously proposed GAT-ST and uses a raw waveform. This captures discriminative cues in both the spectral and temporal domains. The proposed RawGAT-ST model uses a one-dimensional convolution layer to ingest raw audio. This end-to-end architecture uses feature representation learning and a GAT that learns the relationships between cues at different sub-band and temporal intervals. This model uses the raw waveforms from the ASVspoof 2019 dataset and achieves an EER of 1.06% and a min-tDCF of 0.0335. However, this model is only tested on one dataset. Yang and Das, in [107], propose a new feature that can capture discriminative information between natural and spoofed speech. The proposed feature, eCQCC-STSSI, uses CQT to transform speech from the time domain to the frequency domain. A DCT is used to decorrelate the feature dimensions and concentrate the energy of the logarithmic octave power spectrum and logarithmic linear power spectrum, respectively. The two DCT outputs are concatenated to form the eCQCC feature vectors. Even though this feature was tested on both ASVspoof2015 and ASVspoof2017, it has not been tested on other datasets and may fail to properly extract features. The authors of [50] concatenate four sub-features on the ASVspoof2019 dataset. These features, Short-Term Spectral Statistics Information (STSSI), Octave-band Principal Information (OPI), Full-band Principal Information (FPI), and Magnitude-Phase Energy Information (MPEI), are fused to generate the delta acceleration coefficients as features for spoofing detection such as CQSPIC, CQEPIC, and CESPIC. The fusion of the features achieves an EER of 7.63% and a min-TDCF of 0.178, but has been tested in many scenarios. In [108], the authors propose a heuristic feature extraction method based on Multi Level Transform (MLT), which extracts valuable information from the octave power spectrum for spoofing attack detection. It relies on MLT to extract relevant information from previous DCT results. The authors apply it to the ASVspoof2015 and ASVspoof2017v2

datasets, and achieve an EER of 14.45. However, it is only tested on two domains. This paper, [74], compares genuine and spoofed speech across different phonemes and shows that specific phonemes (fricatives, nasals, stops, and pauses) are more informative in the detection of replay attacks. The paper creates four different fusion scoring methods to incorporate phonetic information using phoneme-specific models. This method is tested on the ASVspoof2017 V2 dataset, and achieves an EER of 6.18%. However, it has to be fused with a phenome-independent model for the best results. The work of [88], which also looks at sound characteristics, creates Voice-Pop, which identifies a live user by detecting the pop noise naturally incurred by a user breathing while speaking close to the microphone. The authors use their own dataset with GFCC features and achieve an EER of 5.4%. However, this method was only tested in one domain. A few solutions use lightweight deep learning systems to detect replay spoofing. A lightweight CNN [47], based on Maximal Feature-Map (MFM) activation, is used to detect replay attacks in [97]. MFM is able to minimize the dimensionality by using the most relevant features for classification. This method, [97], is extended in [5] to investigate the efficacy of angular margin-based softmax activation in training a light CNN for cloning and replay spoof detection.

### D. End to End Countermeasures

End-to-end techniques have been shown to be effective in a variety of domains. However, relatively little work has been done in the domain of voice spoofing countermeasures to date. In [17], an end-to-end system was proposed to detect replay attacks against ASVs. This system used a combination of a CNN, LSTM, and DNN to take in raw audio as input and classify the audio into genuine, synthetic/cloned, or replay. This model was evaluated on the ASVspoof 2015 and BTAS 2016 datasets and achieved a half total error rate (HTER) of 1.56%, compared to the GMM-PLP-39 baseline model of 2.96%. This system may be enhanced in the feature extraction stage by employing better time and frequency filtering networks. Recently, in [32], the author proposes a novel heterogeneous stacking graph attention layer that models artifacts spanning heterogeneous temporal and spectral domains with a heterogeneous attention mechanism and a stack node. In concert with a new max graph operation that involves a competitive mechanism and an extended readout scheme, their approach, named ASSIST, achieves an EER of 0.83% and a min-tDCF of 0.0275 on ASVspoof2021. However, the proposed method is only tested on one dataset.

In [77], a graph attention network (GATs) is proposed that works by applying a self-attention mechanism to GNNs and modeling graph-structured data. Each node in the graph is weighted according to its relevance to other nodes. GATs can be used to model a specific sub-band or temporal segment using high-level representations extracted from deep residual networks. It achieves a min-tDCF of 0.0089 on the ASVspoof2019 dataset. However, GAT operates on filter bank outputs, not on the waveform. It is only tested on one dataset.

**TABLE V** Cross corpus evaluations of existing methods

| Author | Corpus | Cross-Corpus |
|--------|--------|--------------|
| [3] | ASVspoof 2019 LA | No |
| [24] | ASVspoof 2019 LA | No |
| [31] | chsc2011 & GTZAN | No |
| [44] | NIST 2001 & 2006 & SRE | No |
| [73] | In-house dataset | No |
| [23] | ASVspoof 2017 v1.0 | No |
| [39] | ASVspoof 2017 | No |
| [60] | ASVspoof 2017 | No |
| [104] | ASVspoof 2015 | No |
| [59] | SPINE2000 | No |
| [67] | NIST 2010 evaluation | No |
| [72] | ASVspoof 2017 dataset | No |
| [22] | ASV spoof 2017 | No |
| [100] | Part I-RSR2015 corpus | No |
| [106] | ASVspoof 2017 v2.0 | No |
| [8] | ASVspoof 2017 | No |
| [55] | VSDC & ASVspoof 2019 | No |
| [13] | ASVspoof 2017, 2019 PA | No |
| [12] | ASVspoof 2017 | No |
| [83] | ASVspoof 2015, 2017 & RedDots | Yes |
| [109] | ASVspoof 2017 V2.0 | No |
| [95] | ASVspoof 2019 LA | No |
| [96] | ASVspoof 2019 LA | No |
| [37] | ASVspoof 2015,2017 | No |
| [90] | ASVspoof 2017 | No |
| [89] | ASVspoof 2019 | No |
| [6] | ASVspoof 2017,2019 & ReMASC | Yes |
| [17] | BTAS2016 dataset | No |
| [57] | ASVspoof 2019 | No |
| [81] | ASVspoof 2017 | No |
| [30] | ASVspoof 2017 | No |
| [71] | ASVspoof 2015 | No |
| [62] | VoxForge | No |
| [80] | ASVspoof 2017 | No |
| [25] | BTAS2016 & ASVspoof2017 | No |
| [45] | ASVspoof 2017 | No |

### E. Discussion

The majority of these papers do not test across multiple datasets to establish the generalizability of the solution, perhaps due to the lack of publicly available datasets at the time of publication. However, there are numerous datasets available for testing right now, and a detailed analysis of cross-corpus evaluation of the solutions is presented in Table V. This analysis shows that there are only two articles [83], [6] that have performed the cross-corpus evaluation to date. However, there are some research gaps in the work of these two papers, and they do not perform a thorough analysis of cutting-edge techniques. Some of these papers use raw wave forms, which can save time in feature extraction but need to be verified under controlled conditions. Also, the depth of each model varies, which can increase accuracy on test sets but take longer to train. In terms of deployment, ideally models need to be lightweight enough to be run on smart speakers. New machine learning techniques for the audio domain can be adapted from the video domain, which is generally more developed. In a world that changes frequently, the models need to be able to be adapted to new challenges, and some of these models are designed for a specific task. New models need to be versatile enough for various conditions. Additionally, these models need to be unbiased, an aspect that has not been tested.

### F. Unified voice spoofing countermeasures

By and large, existing countermeasures for voice spoofing attacks exclusively target one type of attack (e.g., voice replay or speech synthesis), and only recently few works have been reported that utilize a unified approach, capable of dealing with multiple types of voice spoofing attack. In [29], the authors introduce a novel ATCoP feature descriptor for the detection of voice presentation attacks, capable of recognizing both LA and PA attacks. Although a good all-around solution, this unified anti-spoofing technique is the most effective at detecting single-and multi-order replay attacks. The experimental results show that the presented approach is effective when tested on four distinct datasets, with either replay or clone forgeries. Although this method is evaluated against several datasets, the success of the given ATCoP descriptors is not reported for DF attacks. In another work [28], the authors recognized a novel hybrid voice spoofing attack, a cloned replay, which may also be used to spoof an ASV system. A cloned replay attack duplicates the audio of a target speaker and transmits it to an ASV system. The author establishes the basis for a spoofing countermeasure capable of detecting multi-order replay, cloning, and cloned-replay attacks through the use of the proposed ATP-GTCC features. This approach is capable of identifying state-of-the-art (SOTA) voice spoofing attacks with a unified solution.

Rostami and Homayounpour et al. [68] proposed an effective attention branch network (EABN) for detecting LA and PA attacks. The provided technique [68] achieves EERs of 0.86% and 1.89%, and min-TDCFs of 0.02% and 0.50%, against PA and LA attacks, respectively. Despite outperforming SOTA approaches on LA and PA attacks, [68] was not tested on deepfake attacks. In contrast, in [46], SENet and ResNet were combined with statistical pooling to handle anti-spoofing with deeper and faster-trained DNNs. It consisted of a SENet, a median-standard ResNet, a dilated ResNet, and an attentive filtering network. A GNN back end classifier was implemented using CQCC and LFCC features. It obtained an EER of 0.59% for a PA and 6.7% against LA attacks. Because the data was obtained directly from the ASVspoof2019 dataset with no data augmentation or cross-corpus dataset, the proposed method's performance was not checked against DF attacks and may decline in a real-world scenario. In [11] the authors presented Emphasized Channel Attention, Propagation, and Aggregation Time-Delayed Neural Networks (ECAPA-TDNN) as their primary model. The authors' intention was to tackle the issue of channel variability by employing an acoustic simulator in order to enhance the original datasets with transmission codecs, compression codecs, and convolutional impulse responses. The presented method attained an EER of 5.46% and a min-tDCF of 0.3094 in the 2021 LA task and an EER of 20.33% on the DF task. Although the presented work prevents several types of attacks, it needs testing and reporting against PA attacks. To the best of our knowledge, no countermeasures have been reported for all four types of voice spoofing attacks, i.e., LA, PA, DF, and cloned replay attacks. However, the community's attention has lately switched to spoofing-aware ASV systems, as mentioned

in the section below. New integrated spoofing technologies have been presented that conducted speech spoofing and automated speaker verification simultaneously, in addition to unified techniques for cutting-edge spoofing attacks.

### G. Integrated spoofing-aware ASV systems

The research community for secure voice-enabled systems is currently focused on integrating research efforts on speaker verification and anti-spoofing. Unlike existing anti-spoofing systems, which focus on independently streamlined spoofing detection, speaker verification may also be embedded in order to build a integrated system. A solution is considered as an Integrated spoofing-aware ASV systems if it is an integrated and optimized system that protects against spoofing attacks while also performing speaker verification. The first paper in this regard is [35], in which the authors propose the spoofing-aware speaker verification (SASV) challenge, which integrates speaker verification and anti-spoofing. In this challenge, the organizers encourage the development of integrated SASV systems that use new metrics to evaluate joint model performance by releasing official protocols and baseline models. The authors extend speaker verification by including spoofed trials in addition to the standard set of target and impostor trials. Unlike the existing ASVspoof challenge, which focuses on separate spoofing detection and speaker verification systems, SASV aims to develop jointly optimized secure ASV solutions. Open-source, pre-trained spoofing detection and speaker verification models are used in two baseline SASV solutions. Participants have free access to both models and baselines, which can be used to develop back end fusion approaches or end-to-end solutions. The top performing system reduces the equal error rate of a conventional speaker verification system from 23.83% to 0.13% when tested with target, bonafide non-target, and spoofed non-target trials. The SASV challenge results demonstrate the dependability of today's cutting-edge approaches to spoofing detection and speaker verification. In this paper, [111], Zhang et al., develop a probabilistic framework for combining the ASV and CM subsystem scores. In addition to the probabilistic framework, the authors propose direct inference and fine-tuning strategies, based on the framework, to predict the SASV score. Surprisingly, these strategies reduced the SASV EER of the baseline to 1.53% in the official SASV challenge evaluation trials. The author validates the efficacy of the proposed modules through ablation studies and provides insights through score distribution analysis.

The author of [51], expresses concern about the improvement of spoofing robustness in automatic speaker verification (ASV) systems in the absence of a separate countermeasure module. To address this issue, the ASVspoof 2019 baseline model is used, in accordance with the back-end machine learning classifier's probabilistic linear discriminant analysis (PLDA). Three unsupervised domain adaptation techniques are used to optimize the back-end using audio data from the ASVspoof 2019 dataset training partition. The results show significant improvement in both the logical and physical access scenarios, particularly in the latter, where the system is attacked by replayed audio, with maximum relative improvements of 36.1% and 5.3% in bonafide and spoofed cases,

respectively. However, it is observed that absolute error rates on spoof trials remain too high on spoofing attacks. This demonstrates the challenge of making a conventional speaker embedding extractor with a PLDA back-end work on a mix of bonafide and spoofed data.

Over the last few years, research has improved the performance of ASVs and countermeasure systems, resulting in low EERs for each system. However, research on the joint optimization of both systems is still in its early stages. This paper, [82], proposes a Spoof-Aggregated-SASV (SA-SASV), an ensemble-free, end-to-end solution for developing an SASV system with multi-task classifiers. An SA-SASV system is further optimized by multiple losses and more flexible training set requirements. The proposed system is trained using the ASVspoof2019-LA dataset. SA-SASV EER results show that training in complete automatic speaker verification and countermeasure datasets can improve model performance even further. The results show that the SA-SASV feature space outperforms previously published approaches in terms of distinguishing spoof attacks and speaker identification. Furthermore, the SA-SASV EER is reduced from 6.05%, produced by previous state-of-the-art approaches, to 4.86% when no ensemble strategy is used. The article argues that a larger set of data and distinct encoders will further improve the EER.

## V. Countermeasures

This section provides a detailed discussion of the selected countermeasures used for comparative analysis of the voice spoofing countermeasures. For this comparative analysis, we have selected a variety of audio features, including the baseline features of the ASVspoof challenges, which are employed in the voice spoofing detection literature. However, it should be noted that the individual features are compared to one another, excluding feature fusion, to investigate the effectiveness of the features themselves and not their complementary interactions. Each feature was carefully crafted to offer as much coverage as possible while maintaining a relatively small comparative list, and each has its own distinct approach to spoofing detection. To fairly choose this list of auditory features, we conducted a survey of the literature, from 2015 to the present to determine the 14 most effective and commonly used feature descriptors. We employ the frequently used machine learning-based classifiers, i.e., GMM, SVM, and evaluate the performance of the features on the recently released deep learning CNN and CNN-GRU classifiers. The next section discusses the extraction of the SOTA countermeasures.

### A. Constant Q cepstral coefficient (CQCC)) [83]

The ASVspoof 2019 challenge defined two baseline systems; one used the Constant Q cepstral coefficient (CQCC) while the second employed linear frequency cepstral coefficients (LFCC) [83] with a GMM as the back-end classifier. CQCC features use the constant q transform (CQT) that allows the capture of variable spectral resolution, dependent on frequency, i.e., increased temporal resolution at higher frequencies and lower at lower frequencies. This

variable resolution analysis allows for speech characteristics to be better captured at differing frequencies. A CQCC extraction was performed by first applying CQT to the windowed audio sample, and then a fast fourier transform (FFT) was used to obtain the power spectrum of the obtained windows, which were then log-scaled. After log scaling, uniform re-sampling was performed, and finally a discrete cosine transform (DCT) was used to obtain the final cepstral coefficients. The extraction of CQCC features was performed using the following method:

$$X^{CQ}(k,n) = \sum_{j=n-\lfloor N_k/2 \rfloor}^{n+\lfloor N_k/2 \rfloor} x(j)a_k(j-n+N_k/2) \quad (1)$$

$$a_k(n) = \frac{1}{C}\left(\frac{n}{N_k}\right)exp[i(2\pi n \frac{f_k}{f_s} + \Phi_k)] \quad (2)$$

$$C = \sum_{l=-\lfloor N_k/2 \rfloor}^{N_k/2} w\left(\frac{l+N_k/2}{N_k}\right) \quad (3)$$

$$f_k = f_1 2^{\frac{k-1}{B}} \quad (4)$$

$$Q = \frac{f_k}{f_{k+1} - f_k} = (2^{1/B} - 1)^{-1} \quad (5)$$

$$N_k = \frac{f_s}{f_k}Q \quad (6)$$

$$CQCC(p) = \sum_{l=1}^{L} log|X^{CQ}(l)|^2 cos[\frac{p(l-\frac{1}{2})\pi}{L}] \quad (7)$$

where $a(n)$ represents a complex conjugate of $a(n)$, $N_k$ denotes a variable window lengths, w(t) is a window function (i.e. the Hanning Function). $f_k$ is the central frequency, f1 refers to the center frequency of the lowest bin, and B determines the number of bins, or octaves.

## B. Mel-Frequency cepstral coefficients (MFCC) [58]

The Mel-Frequency Cepstral Coefficients (MFCC) [58] feature is well known and has been used successfully in various forms of audio recognition. To obtain MFCC features, the audio sample is divided into a number of equal windows. Following that, an FFT is performed to acquire the spectrum, which is then passed through a Mel-scaled filterbank made of triangle bandpass filters. Log-scaling is then performed, and the cepstral coefficients are obtained using a DCT.

MFCCs are computed as follows:

$$MFCC(q) = \sum_{m=1}^{M} log[MF(m)]cos[\frac{q(m-\frac{1}{2})\pi}{M}] \quad (8)$$

$$MF(m) = \sum_{k=1}^{K} |X^{DFT}(k)|^2 H_m(k) \quad (9)$$

where k refers to the discrete fourier transform (dft) index and H(k) denotes the triangular weighting-shaped function for the m-th Mel-scaled band-pass filter.

## C. Inverse Mel-Frequency cepstral coefficient (IMFCC) [9]

One variation of MFCC, the inverse Mel-Frequency cepstral coefficients [9], has also been employed for voice spoofing. These features are capable of capturing important traits of the audio signals by analyzing the high frequency components by using an inverted Mel-scaled filter bank made up of triangular band pass filters. IMFCC is obtained as follows:

$$IMF(m) = \sum_{k=1}^{K} |X^{DFT}(k)|^2 H_m(k) \quad (10)$$

$$IMFCC(q) = \sum_{m=1}^{M} log[IMF(m)]cos[\frac{q(m-\frac{1}{2})\pi}{M}] \quad (11)$$

where k and H(k) refer to the index and triangular weighting-shaped function for the m-th Inverted Mel-scaled band pass filter.

## D. Linear Frequency Cepstral Coefficient (LFCC) [113]

Another common variant of MFCC is the Linear Frequency Cepstral Coefficient (LFCC) [113], which is derived via the same process as MFCC but uses a linear frequency filter bank in place of the Mel filter bank. In the linear frequency filter bank, the bandwidth of the triangular filter remains constant, in contrast to the Mel filter bank, which changes slope over the frequency range. Because of this, LFCC has been proven to be more robust than MFCC in detecting artifacts at higher frequencies[31].

$$f(m) = \left(\frac{N}{F_s}\right)(f_1 + m\frac{f_h - f_l}{M+1}) \quad (12)$$

where f(m) are the frontier points of the linear filterbank.

## E. Magnitude-based Spectral Root Cepstral Coefficients (MSRCC) and Phase-based Spectral Root Cepstral Coefficients (PSRCC) [80]

Magnitude-based Spectral Root Cepstral Coefficients (MSRCC) and Phase-based Spectral Root Cepstral Coefficients (PSRCC) are spectral root features [80] that capture both the magnitude and phase information, respectively, to produce distinct traits in bonafide and spoofed audio. These features are computed as follows:

$$MSRCC(q) = \sum_{m=1}^{M} (MFM(m))^\gamma cos[\frac{q(m-\frac{1}{2})\pi}{M}] \quad (13)$$

$$MFM(m) = \sum_{k=1}^{K} |S(k)|^2 H_m(k) \quad (14)$$

$$PSRCC(q) = \sum_{m=1}^{M} (MFP(m))^\gamma cos[\frac{q(m-\frac{1}{2})\pi}{M}] \quad (15)$$

$$MFP(m) = \sum_{k=1}^{K} \angle S(k)H_m(k) \tag{16}$$

where $S(k)$ is the k-point of the DFT of signal $s(n)$ and $H_m(k)$ is the triangular weighting-shaped function for the m-th Mel-scaled bandpass filter.

### F. Spectral Centroid Magnitude Coefficients (SCMC) and Spectral Centroid Frequency Coefficients (SCFC) [44]

Spectral Centroid Magnitude Coefficients (SCMC) and Spectral Centroid Frequency Coefficients (SCFC) are the two-dimensional representations of sub-band energy in the speech spectrum [44]. These two-dimensional features include the centroid extraction of magnitude and frequency-based features from the speech spectrum. Spectral centroid features were previously employed for speaker recognition, and sub-band spectral features have demonstrated significant efficacy in noisy speech recognition. Centroid features carry formant-related information, which has been shown to be robust to the presence of noise in the vocal sample [44]. Spectral centroid frequency (SCF) is the weighted average frequency for a specific sub-band, where the weights are the normalized energy of each frequency component of the sub-band. Spectral centroid magnitude (SCM) is the weighted average magnitude for a given sub-band. These coefficients are derived as follows:

$$SCM_t^i = \frac{\sum_{k=1}^{M/2+1} f(k)S_t(k)w_i(k)}{\sum_{k=1}^{M/2+1} S_t(k)w_i(k)} \tag{17}$$

$$SCF_t^i = \frac{\sum_{k=1}^{M/2+1} f(k)S_t(k)w_i(k)}{\sum_{k=1}^{M/2+1} f(k)w_i(k)} \tag{18}$$

where $S_t(k)$ and $f(k)$ are the power spectral magnitude of the t-th frame, and normalized frequency corresponding to the k-th frequency component, respectively. SCFCs are extracted directly using linear rectangular filter banks, whereas SCM is passed through log scaling and DCT to obtain SCMC features.

### G. All-Pole Group Delay Function (APGDF) [67]

Unlike the features derived from the magnitude spectrum, e.g., MFCCs, the All-Pole Group Delay Function (APGDF) is extracted from the phase spectrum, which is largely ignored by the commonly used cepstral coefficients, e.g., GTCC, LPCC, and CQCC [67]. APGDF provides an efficient way to utilize the information from the voice signal's phase spectrum. Group delay functions employ parametric all-pole models instead of spectral transforms. APGDFs enhance recognition accuracy and give comparable results to traditional magnitude-based MFCC features by integrating all-pole analysis, filters, and group delay functions. Moreover, APGDF features are computationally less complex than the recently cited MGDF features [59]. APGDF features are extracted as follows:

$$APGDF(i) = \frac{A_P(G)D_P(G) + A_I(G)D_I(G)}{|A(P)|^2} \tag{19}$$

$$A_P = |A_P|\sigma^{i\vartheta(P)} \tag{20}$$

Where $A_P(G)$ and $D_P(G)$ are Fourier transforms, $|A_P|$ refers to the magnitude spectrum and $i\vartheta(P)$ represents the phase spectrum. Finally, a DCT is used to convert the group delay function into cepstral coefficients [67].

### H. Rectangular Filter Cepstral Coefficient (RFCC)

Another tested feature is the Rectangular Filter Cepstral Coefficient (RFCC). Similar to MFCC features, RFCCs are extracted using a rectangular window spaced on a linear scale instead of triangular filters on the mel-scale. RFCC features were initially proposed for use in automatic speaker recognition systems operating under noisy conditions.

### I. Linear prediction cepstral coefficients (LPCC) [91]

Linear Prediction Cepstral Coefficient (LPCC) features have been historically used to capture the emotion information found in speech. Unlike other cepstral features, LPCCs use linear prediction analysis prior to cepstral analysis. The audio signal first undergoes band limiting, then a high-emphasis filter is applied. After high-emphasis, a 30ms hamming window is applied to the windows at a 10ms interval. The next step is to calculate the first through tenth linear predictor coefficients through the use of auto-correlation. Next, these linear predictor coefficients are transformed into the cepstral coefficients by using the following equations:

$$\begin{aligned} c_1 &= a_1 \\ c_n &= \sum_{k=1}^{n-1}(1 - k/n)a_k c_{n-k} + a_n \quad 1 < n \le p \end{aligned} \tag{21}$$

$$h(t) = \sum_{k=1}^{N} A_k cos(\phi_k(t)) \quad \phi_k = 2\pi k f_0 t + \theta_k \tag{22}$$

where $c_i$ and $a_i$ are the i-th cepstral and linear predictor coefficients, respectively.

### J. Gammatone Cepstral Coefficients (GTCC) [85]

Gammatone Cepstral Coefficients (GTCCs) [85] are similar to MFCC features, but are much more robust to noise than MFCC, offering superior performance in classification tasks. A gammatone filter provides more frequency components in the low-frequency range with narrow bandwidth while providing fewer frequency components in the high-frequency range with wider bandwidth, effectively revealing the spectral information. GTCC features are computed by applying an FFT to each frame in order to get the spectrum. The spectra of each window are then passed through a gammatone filter bank to obtain the energy at each subband. A logarithmic function is then applied to these energies, and finally, DCT is applied to obtain the GTCCs. 13 to 20 coefficients are usually extracted and are considered optimal for audio signal analysis. GTCCs are extracted using the following equation:

$$GTCC_m = \sqrt{\frac{2}{N}} \sum_{n=1}^{N} log(X_n)cos[\frac{\pi n}{N}(m - \frac{1}{2})] \quad 1 \le m \le M \tag{23}$$
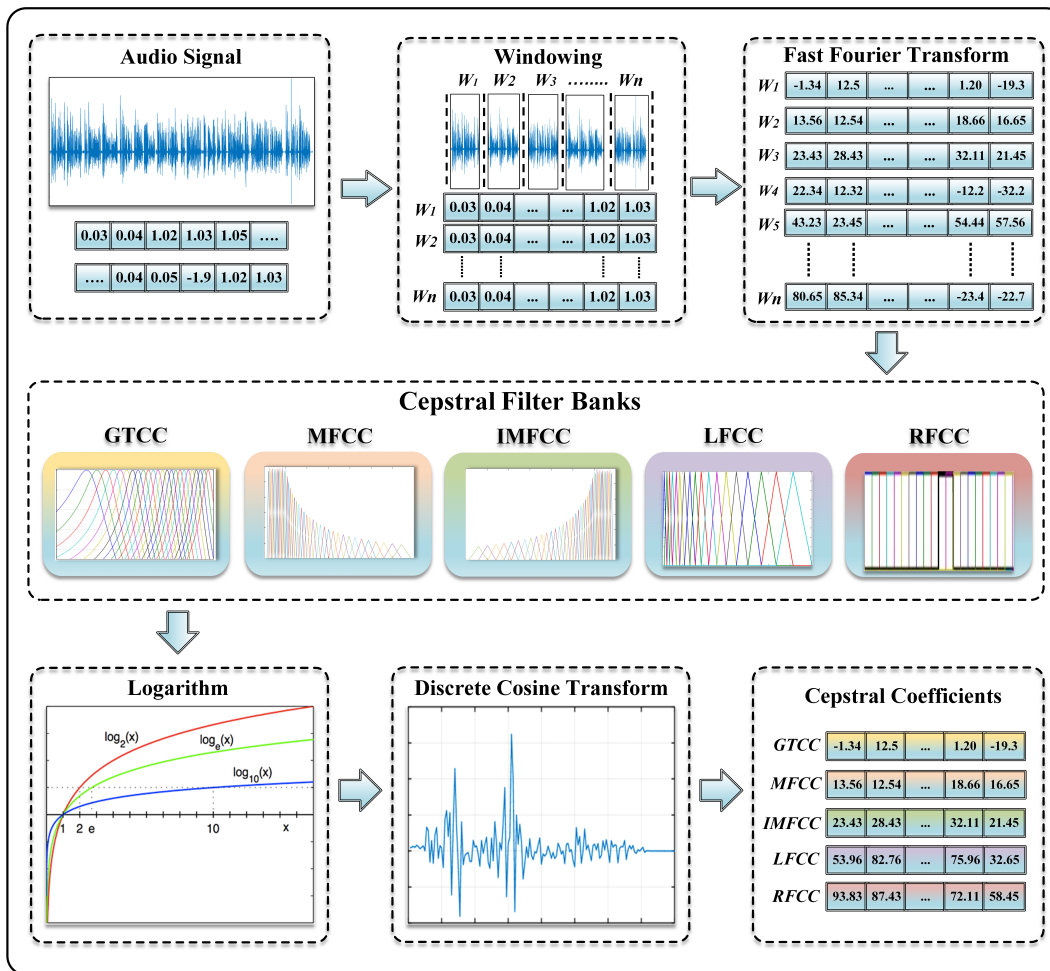
**Fig. 5:** Extraction of GTCC, MFCC, IMFCC, LFCC and RFCC Features.

where $X_n$ is the signal energy of the nth spectral band, N is the number of gammatone filters used, and M is the number of GTCCs to be extracted.

### K. Relative Phase Shift (RPS) [73]

Relative Phase Shift (RPS) features [73] are a representation of harmonic phase information. A harmonic analysis model represents each frame of a signal as a sum of sinusoidal harmonically related to pitch and fundamental frequency. The RPS representation calculates the difference between the instantaneous phase of every harmonic and fundamental component at a fixed point over the period of the signal. RPS features reveal a structured pattern in the phase information of voiced segments, and can be derived as follows:

$$h(t) = \sum_{k=1}^{N} A_k cos(\phi_k(t)) \quad \phi_k = 2\pi k f_0 t + \theta_k \quad (24)$$

A graphical representation for the extraction of features is presented in Figure 5 and 6. The next section contains the experimental details, data sources and the system information used to perform the experimentation and the comparative analysis of the features.

## VI. Publicly Available Datasets

This section of the paper discusses and identifies the various datasets used in cutting-edge ASV systems. Early stages of voice spoofing detection research involved speech and speaker recognition databases, i.e., YOHO [43], NIST [1], and WSJ [18]. However, to accurately account for research progress, there was a dire need for a common dataset as well as a performance metric to evaluate spoofing countermeasures. Consequently, this need was discussed and addressed at the INTERSPEECH 2013 special session on spoofing and ASV countermeasures. This special session motivated the research community to organize the first Automatic Speaker Verification Spoofing and Countermeasures Challenge, ASVspoof, in 2015, which took place at INTERSPEECH that year. The dataset released for this challenge included two types of spoofing attacks: synthetic speech (SS) and voice conversion (VC). In the years following, three additional ASVspoof challenges were organized: ASVspoof 2017, ASVspoof 2019, and ASVspoof 2021, each with publicly available datasets for download. There are several common publicly available datasets which are used by voice PAD researchers, i.e., ASVspoof, AVspoof, ReMASC, Spoofing and anti-spoofing (SAS) corpus, RedDots, Vox Celeb, voicePA, and BioCPqDPA, among others. We made a concerted effort to cover all of the existing datasets
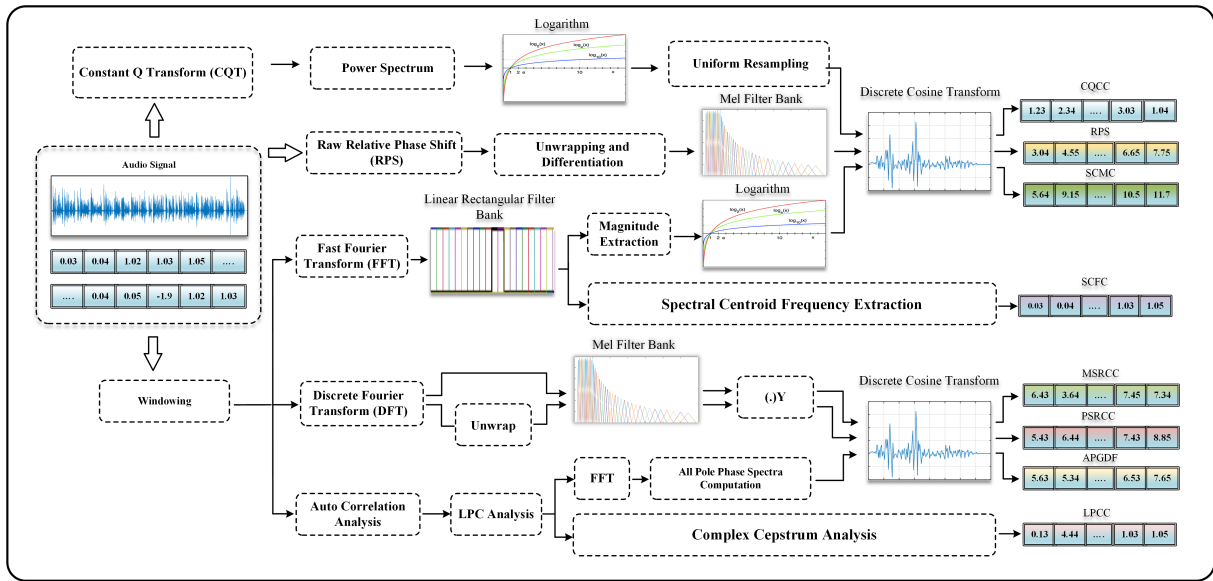
**Fig. 6:** Extraction of CQCC, RPS, SCMC, SCFC, MSRCC, PSRCC and APGDF Features.

(2015–2022) used for spoof detection and countermeasure development. Details of the publicly available datasets which address spoofing attacks are presented in Figure 7.

### A. Spoofing and Antispoofing (SAS) corpus [101]

The Spoofing and Antispoofing (SAS) corpus [101] contains a wide range of spoofed speech samples generated using nine different approaches, two of which are speech synthesis, and the other seven are voice conversion. This database [101] contains two protocols: one for testing the ASV system and another for generating spoofed speech sounds. Consequently, they enable the speech synthesis community to create spoofing materials attentively while remaining unaware of speaker verification spoofing and anti-spoofing. Periods of silence not found in natural speech were removed from the samples, resulting in a more realistic Speech Synthesis (SS) and Voice Conversion (VC) spoof corpus (Wu et al., 2015a, 2015b, 2015c). The SAS corpus contains speech produced using various spoofing methods, represented in 300,000 samples of each type. Without using a countermeasure, ASV systems were extremely vulnerable to spoofing attacks in SAS [101].

### B. RedDots [48], [40]

The RedDots project [48] was launched as a follow-up to a special session at INTERSPEECH 2014, with collaboration from multiple sites. Its goal is to collect speech data through mobile crowd sourcing, which allows for a larger population and greater diversity. On January 29, 2015, the project was launched. At the time of our investigation, the project had 89 speakers from 21 countries, 72 men and 17 women, for a total of 875 complete sessions. The purpose of this special session is to bring together research efforts to investigate potential approaches and gain a better understanding of speaker/channel phonetic variability. The RedDots dataset [48] is intended to have a higher number of recordings per session, with shorter

sessions spent on each. One goal is to collect 52 sessions per speaker, one session per week, for a year. Because of this, each session is limited to two minutes, with a total of 24 sentences (10 common, 10 unique, 2 free choices, and 2 free texts) for each session. This dataset [48] includes a diverse set of inter- and intra-speaker types. Following that, the replayed RedDots database [40] was created by re-recording the original corpus utterances under different environmental conditions. Both of these databases assist in the development of replay-resistant ASV systems because the original RedDots provided genuine utterances. An overview of the RedDots database [48] is presented in Table VII.

### C. AVspoof [42]

The AVspoof dataset is designed to assist ASV systems in the development of anti-spoofing techniques. This database was used in the BTAS 2016 [42], [18] challenge. The AVspoof database includes replay spoofing attacks in addition to synthetic speech and VC spoofing attacks. The replay attacks are generated by various recording devices. The SS attacks were generated by the following techniques, i.e., Hidden Markov Model (HMM) and Festvox, which account for the vast majority of the VC attacks. These sessions are recorded by participants in a variety of environments and with a variety of recording devices. Speakers are instructed to read out sentences, phrases, and speak freely about any topic for 3 to 10 minutes. To make competition more difficult, "unknown" attacks are included in the test set [42]. The organizers of the challenge provide a baseline evaluation system based on the open source Bob toolbox [42]. In the baseline system, simple spectrogram-based ratios serve as features, and logistic regression is used as a pattern classifier [42]. The statistics of the database are summarized in Table VI.

**TABLE VI** A summary of AVspoof Database.

| Subset | Genuine | Spoof (PA) | Spoof (LA) |
|---|---|---|---|
| Training | 4973 | 38580 | 17890 |
| Development | 4995 | 38580 | 17890 |
| Evaluation | 5576 | 43320 | 20060 |

### D. VoxCeleb [61]

VoxCeleb [61] is an audio-visual dataset comprised of short clips of human voices extracted from YouTube interview videos. Each segment lasts at least 3 seconds. VoxCeleb features speech from people of various ethnicities, accents, professional backgrounds, and ages. 61% of speakers are male and 39% are female. The data was collected randomly, with ambient noise, laughter, overlapping speech, pose deviation, and a variety of lighting conditions. This dataset is available in two versions: VoxCeleb1 [61] and VoxCeleb2 [14]. Each has audio files, face clips, metadata about speakers, and so on, in the training and testing sets. The Finnish-language sets assist ASV systems in the detection of mimicry attacks [86]. The details of the dataset [61] are shown in table VII.

### E. voicePA [41]

The voicePA dataset was created with the assistance of the AVspoof dataset. Its bonafide data is a subset of the AVspoof dataset's genuine data, uttered by 44 speakers in four recording sessions held in various settings [89]. These sessions were recorded using high-quality microphones from a laptop, a Samsung S3, and an Apple smartphone 3GS. Spoofed data consists of 24 different types of presentation attacks that are captured using five devices in three different environments. These spoof utterances are based on real-world data. The original dataset's SS and VC spoofed audio samples are also replayed in the voicePA dataset [41]. Table VII contains detailed information about the voicePA [41] dataset.

### F. BioCPqD-PA [41]

The Portuguese language BioCPqD-PA dataset [41] was collected by recording 222 people in a variety of environmental conditions. The dataset was comprised of 27,253 authentic recordings and 391,687 samples, which had been subjected to a presentation attack. In the presentation audio data, one laptop was used with 24 different setups consisting of 8 loudspeakers and 3 microphones, while another single laptop was used to capture real data. This dataset [41] was divided into three parts: training, development, and evaluation. Each part was collected using a different set of microphones and loudspeakers, and each set was comprised of the audio recording from all of the participating speakers. An overview of the BioCPqD-PA [41] dataset is presented in Table VII.

**TABLE VII** A summary of the existing databases.

| Dataset | Language | Attacks | Speakers | Utterances |
|---|---|---|---|---|
| RedDots | 5 | Replay | 3750 | 5000 |
| VoxCeleb | 6+ | Mimicry | 1251 (V1), 7000+ (V2) | 100,000 (V1), 1,000,000 (V2) |
| VoicePA | - | SS, VC | 44 | - |
| BioCPqD-PA | 1 | 222 | 3750 | 418,940 |

**TABLE VIII** A summary of ASVspoof Challenge 2015 database.

| | Speaker | | Utterances | |
|---|---|---|---|---|
| Subset | Male | Female | Genuine | Spoofed |
| Training | 10 | 15 | 3750 | 12625 |
| Development | 15 | 20 | 3497 | 49875 |
| Evaluation | 20 | 26 | 9404 | 193404 |

### G. ASVspoof 2015 [102]

The ASVspoof 2015 Challenge database is the first significant release for research into spoofing and countermeasures [102]. For LA attacks, the database contains natural and spoofed speech generated by speech synthesis and VC. There are no noticeable effects from the channel or background noise. The database is split into three sections: training, development, and evaluation. The evaluation subset is made up of both known and unknown attacks. The known attacks contain the very same five algorithms that were used to generate the development dataset and are thus referred to as known (S1-S5) attacks. Other spoofing algorithms are also included in unknown (S6-S10) attacks that are directly used in the test data. Table VIII contains a detailed description of the challenge database.

### H. ASVspoof 2017 [38]

The ASVspoof 2017 Challenge dataset [38] was constructed on the RedDots dataset [48]. The dataset [38] contains replayed data samples with text-dependent speech. It included the voices of 42 different speakers, recorded using 61 combinations of distinct recording devices, replay devices, and environmental conditions. It was collected over the course of 179 sessions. The original ASVspoof 2017 database [38] contained some inconsistencies, but the issue was resolved in ASVspoof 2017 Version 2.0 [16]. In addition to the corrected data, a more detailed description of recording and playback devices, as well as acoustic environments, was provided [16]. The details of the ASVspoof 2017 [38] dataset are presented in Table IX.

**TABLE IX** A summary of ASVspoof Challenge 2017 database version 2.0

| Subset | Speaker | Genuine | Spoofed |
|---|---|---|---|
| Training | 10 | 1507 | 1507 |
| Development | 8 | 760 | 950 |
| Evaluation | 24 | 1298 | 12008 |

### I. ReMASC [21]

The ReMASC (Realistic Replay Attack Microphone Array Speech Corpus) is a database of speech recordings developed to support research into and the security of voice-controlled systems. The ReMASC database is comprised of authentic and replayed recordings of speech samples, captured in actual circumstances, and utilizes cutting-edge voice assistant development kits. In particular, the ReMASC dataset contains recordings from four systems, each with a different transmitter and receiver, under a range of atmospheric situations, with varying levels of background noise and relative speaker-to-device locations. This is the first database that was specifically

developed to safeguard voice-controlled systems (VCS) from various types of replay attacks in varied contexts. Table X contains the sampling details of the eMASC [21] dataset.

**TABLE X** A summary of ReMASC Database.

| Environment | Subject | Genuine | Replayed |
|---|---|---|---|
| Outdoor | 12 | 960 | 6900 |
| Indoor1 | 23 | 2760 | 23104 |
| Indoor2 | 10 | 1600 | 7824 |
| Vehicle | 10 | 3920 | 7644 |
| Total | 55 | 9240 | 45472 |

### J. ASVspoof 2019 [89]

The ASVspoof 2019 challenge [89] is an extension of the previously held ASVspoof challenges and focuses on countermeasures for all three major attack types, namely, SS, VC, and replay. It is divided into the logical access (LA) and physical access (PA) subsets. LA contains TTS and VC spoof speech samples, and PA has replay-spoof speech. Both of these subsets are further partitioned into training, development, and evaluation subsets. The training subset is generated by 8 males and 12 females; the development subset by 4 males and 6 females, and the evaluation subset by 21 males and 27 females. The spoofed speech signals are generated using one of the two VC and four speech synthesis algorithms. Details of the ASVspoof 2019 [89] dataset are shown in Table XI.

### K. ASVspoof 2021 [15]

ASVspoof 2021 [15] is the fourth offering in the series of spoofing challenges designed to promote spoofing research and the development of countermeasures to protect automatic speaker verification systems from manipulation. ASVspoof 2021 [15] introduces a new task involving deepfake speech detection, in addition to a continued focus on logical and physical access tasks, and has a number of advancements over previous editions. ASVspoof 2021, in particular, is divided into three sub-challenges. The first is a logical access sub-challenge that builds on the 2019 challenge by emphasising robustness to channel variation; the second is a physical access sub-challenge, similar to the 2019 setup, but with recordings made in real-world physical environments. The third is a speech deepfake detection sub-challenge (no ASV). ASVspoof 2021 contains technically difficult data to encourage broad generalization of countermeasures.

**TABLE XI** A summary of ASVspoof Challenge 2019 database.

| | Speaker | | LA Attacks | | PA Attacks | |
|---|---|---|---|---|---|---|
| Subset | Male | Female | Genuine | Spoofed | Genuine | Spoofed |
| Training | 8 | 12 | 2580 | 22800 | 5400 | 48600 |
| Development | 8 | 12 | 2548 | 22296 | 5400 | 24300 |
| Evaluation | - | - | 71747 | | 137457 | |

The logical access (LA) task is similar to that of ASVspoof 2017 and 2019, but this version includes the computation and transmission of text-to-speech (TTS) and voice conversion (VC) attacks. In comparison to ASVspoof 2017, the physical access (PA) task includes genuine and replayed samples but

with a more tightly controlled setup. The new speech deepfake task (DF), similar to the LA task, includes compressed audio. The protocols used in the training and development sections are the same as they were in ASVspoof 2019, and therefore ASVSpoof 2021 does not include development or training subsets. There are new metrics for the evaluation partition, including a slightly revised t-DCF metric for the LA and PA tasks. Whereas the EER is used for the evaluation of DF tasks.

### L. Voice Spoofing Detection Corpus (VSDC) [6]

The VSDC dataset [6] was designed to be a standard dataset comprising both real audio recordings from various surroundings and microphones, as well as spoofed audio files generated through various controlled environments, spoofing realistic scenarios. This data can be analyzed and utilized to develop countermeasures to avoid audio spoofing attacks in the replay category. Since an Internet of Things (IoT) device may be exploited as a point of replay to another device, audio replay attacks have become more significant. This offers a perfect atmosphere for carrying out repeat attacks. This data set was generated with the intention of simulating these attacks in a controlled environment. While the primary goal of this dataset is to detect multi-hop replay attacks, it is not restricted to that. It could be used to investigate typical replay attacks, the influence of different microphones and surroundings on an audio file, or how an individual's vocal range affects the accuracy of a voice control system. The details of the number of sample and development architecture can be found in [6].

## VII. Performance Evaluation Parameters

This section describes the performance evaluation criteria used in this survey, and the specific metrics are defined in the following subsections. After an extensive review of published work, we observed that Equal Error Rate (EER) was the primary criterion used for the performance evaluation of voice spoofing and ASV systems. Of the articles surveyed for this paper, more than three-quarters evaluated their presented work using EER, for instance [26]. A small number of the investigated papers picked accuracy as the sole performance rating criterion, however, the majority of recent work uses multiple rating criteria to evaluate the performance of the presented method, e.g., EER with min-tDCF, EER with Half Total Error Rate (HTER), and False Match Rate (FMR) with False Non-match Rate (FNMR). As not all research works were examined using the same criteria, comparative analysis of the results became problematic.

### A. Equal error rate (EER)

An ASV system either approves or rejects the specified identities. A categorization might be correct or erroneous in four ways. These are True Acceptance (TA), True Rejection (TR), False Acceptance (FA), and False Rejection (FR). TA and TR are favorable outcomes, but FR and FA are damaging to the system. These possibilities are determined by a preset threshold [83]. In the case of FA, a faked speech sample with a score higher than or equal to the preset threshold is allowed.
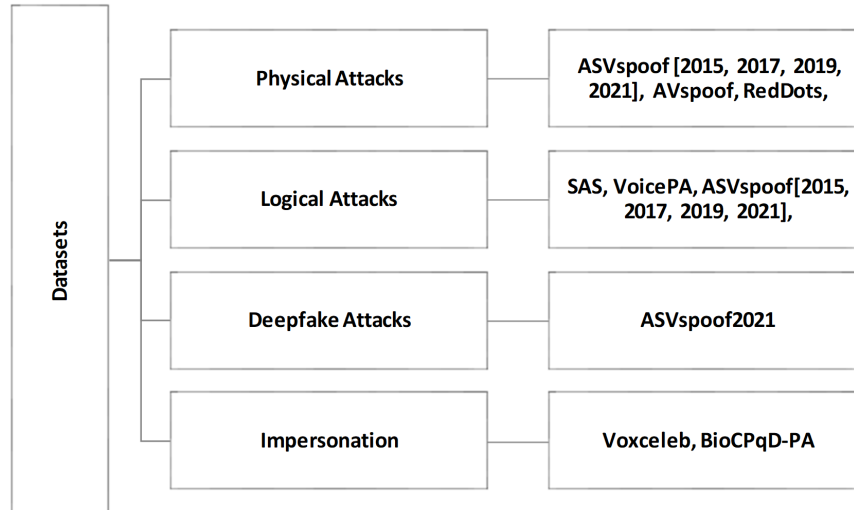
**Fig. 7:** Existing Datasets used for Automatic Speaker Verification Systems.

Conversely, in the case of FR, a speech sample with a score less than the predefined threshold is rejected. To evaluate ASV performance, the Equal Error Rate (EER), the value at which the False Acceptance Rate (FAR) and False Rejection Rate (FRR) become equal, is employed.

$$FAR = \frac{FA}{TA} \tag{25}$$

$$FRR = \frac{FR}{TT} \tag{26}$$

where FAR and FRR denote the false acceptance and false rejection rate, FA and TA represent the number of false acceptance and false rejection, and TA and TT refer to the total number of false and genuine speech samples.

### B. Detection Error Tradeoff (DET) curve

A Detection Error Tradeoff (DET) graph is a graphical representation of errors for classification model systems that compares the false rejection and the false acceptance rate. The x- and y-axes are scaled non-linearly by their standard normal deviations (or simply by logarithmic transformation), resulting in tradeoff curves that are more linear than ROC curves and use the majority of the area to highlight substantial differences in the crucial operating zone. With its effectiveness for binary classification, ASV systems are initially evaluated using the DET curve. It generates EER curves by plotting FAR on the x-axis and FRR on the y-axis.

### C. Half Total Error Rate (HTER)

Spoofed speech with a higher score than the predefined threshold will be misidentified as bonafide samples, whereas bonafide samples with a lower score will be incorrectly classified as spoofed speech samples. Since the two errors, FAR and FRR, are inversely connected, it is essential to depict performance as a function of the threshold. The Half Total Error Rate is one such metric (HTER). HTER's computation is shown below:

$$HTER = \frac{FAR + FRR}{2} \tag{27}$$

### D. Tandem Detection Cost Function (t-DCF)

The ASVspoof challenge series was created to lead anti-spoofing research for automated speaker verification (ASV). The 2015 and 2017 challenge editions featured evaluating spoofing countermeasures (CMs) in isolation from ASV using an equal error rate (EER) metric. While this was a strategic method of evaluation at the time, it had certain flaws. First, when ASV and CMs are coupled, the CM EER is not always a trustworthy predictor of performance. Second, the EER operating point is inappropriate for user authentication systems with a high target user prior but a low spoofing attack prior, such as telephone banking. As a consequence, the community will shift from CM-centric to ASV-centric evaluation using a new tandem detection cost function (t-DCF) measure. It extends the typical DCF used in ASV research to spoofing attack circumstances. The t-DCF metric is made up of six components in two parts: (Part i) false alarm(1-2) and miss costs(3-4) for both systems; and (Part ii) prior probability of target(5) and spoof trials(6) (with an implied third, non-target prior). The results of t-DCF in [37] justify the inclusion of the DCF-based metric in the roadmap for future ASVspoof challenges, as well as additional biometric anti-spoofing evaluations.

## VIII. Comparative analysis of voice spoofing countermeasures

In this section of the paper, we discuss the datasets and classifiers used to differentiate between spoofed and real speech samples. The dataset subsection contains information on speech samples and corpus details that were utilized in training, development, and testing, as well as cross-corpus evaluation of countermeasures. The subsection on feature and classifier configuration describes the back end classification approach used to evaluate the effectiveness of the countermeasures. Next, the experimental analysis's environmental setup is described.

## A. Dataset

For the comparative analysis, we used two datasets to perform the experimentation: VSDC [6] and ASVspoof 2019 [89]. The VSDC dataset consists of a wide variety of microphones and speakers. Although the VSDC dataset contains types of bonafide, first-order, and second-order replay recordings, it is designed to provide a diverse set of recording and playback environments. Specifically, the microphones range from professional grade to cellphones, and the speakers for replay range from professional to local residents. In the VSDC, a total of fifteen (15) speakers, ages 18–60 years old, participated in data collection. Out of fifteen (15) speakers, ten (10) are male, and five (5) are female. Some of the speakers are not native English speakers. Each speaker recorded the original file by repeating a given set of phrases typical of commands given to VCDs. Some of the volunteers recorded these original phrase sets multiple times using different microphones in diverse environments. In total, 198 different 0PR source sets, consisting of at least 9 phrases, were recorded, resulting in 1,926 0PR source phrases being spoken. 1PR and 2PR samples were obtained from 0PR samples. There were 22 different microphones used, with 14 different microphone-speaker configurations for 1PR and 8 different configurations used for 2PR. 1PR and 2PR each include 4,923 samples, giving VSDC a total of 11,772 samples.

In contrast, ASVspoof 2019 is the larger dataset, consisting of 107 speakers. The ASVspoof 2019 dataset includes bonafide and spoof recordings, with LA and PA attacks. The recording environments are more uniform and less diverse than in VSDC, consisting of only hemi-anechoic chambers of varying room sizes. The ASVspoof2019 dataset brings additional diversity in the length of the audio recordings, both bonafide and spoofed, that VSDC does not. The ASVspoof 2019's PA dataset consists of 54,000 samples in the training data set, 29,700 samples in the validation set, and 137,457 samples in the evaluation set. The samples were recorded from 20 individuals, consisting of 12 female and 8 male participants. The recordings utilized a 16-kHz sampling rate and 16-bit resolution, along with having 27 different acoustic configurations. The configurations had variations using 3 room sizes, 3 different levels of reverberation, and 3 different speaker-to-ASV distances. For the replay sample of the data set, 9 different configurations were used, each having 3 different attacker-to-talker distances and 3 different speaker quantities. It should also be noted that configurations for the testing data set differed from those of the training and validation sets. The more detailed description of the dataset is provided in section VI and subsection J.

## B. Feature and classifier Configuration

To test the effectiveness of the state-of-the-art (SOTA) spoofing countermeasures, four different classifiers (two machine learning and two deep learning) comprised of a Gaussian Mixture Model (GMM) [84], a Support Vector Machine (SVM-rbf) [55], a Convolutional Neural Network (CNN) [34], and a CNN-Gated Recurrent Unit (CNN-GRU) [34] were used. The baseline GMM model given by ASVspoof2019 for CQCC and LFCC features was expanded to compute the

**TABLE XII** Configuration details of the comparative countermeasures.

| Features | | |
|---|---|---|
| Feature | Type | Configuration |
| CQCC[83] | Cepstral | No. of FFT bins (per octave)=96 |
| LFCC [113] | Cepstral | No. of Filters=20 |
| MFCC [58] | Cepstral | No. of Filters=14 |
| IMFCC [9] | Cepstral | No. of Filters=20 |
| SSFC [80] | Cepstral | No. of Coeffs.=20 |
| RFCC [73] | Cepstral | No. of Filters=20 |
| PSRCC [80] | Cepstral | No. of Coeffs.=13 |
| MSRCC [80] | Cepstral | No. of Filters=13 |
| LPCC [91] | Cepstral | LP order=20 |
| SCMC [44] | magnitude | No. of Sub-bands=20 |
| SCFC [44] | frequency | No. of Sub-bands=20 |
| GTCC [85] | Cepstral | No. of Coeffs.=14 |
| RPS [73] | Cepstral | No. of Coeffs.=20 |
| APGDF [67] | Phase based | LP Order=20 |

performance of all frontend features. The GMM was trained for ten iterations, and the score of the speech samples was computed using a log-likelihood ratio. EER and min-TDCF were used to assess the effectiveness of each countermeasure with GMM baseline classification. Next, an SVM classifier with an RBF kernel was used to evaluate the performance of the SOTA countermeasures. The GMM and SVM classifiers were chosen because they had the best reported results for the respective datasets. For example, the SVM-rbf classifier [55] achieved the best performance on the VSDC dataset, whereas CQCC-GMM and LFCC-GMM were the baseline classifiers for the ASVspoof2019 competition. In addition, CNN and CNN-GRU [34] were used to determine the performance of the SOTA countermeasures on deep learning-based classifiers. The effectiveness of each countermeasure is described below in the experiment section. Section V also completely describes the extraction of the countermeasure chosen for comparative study. Although we assured the codes' availability for reproducibility, the two-dimensional characteristics of the features was converted to one-dimensional by employing the mean average of the retrieved features. The features are detailed in the code and the ReadMe.txt file. The code and ReadMe.txt file contain details of the features. Further details of the extraction mechanism used for the frontend features are mentioned in Table XII.

## C. Experimental Setup

The setup for the experiments consisted of a pipeline of feature extraction and then evaluation of the extracted features on the ML and DL-based classifiers. The feature extraction and classifiers were run on Oakland University's Matilda High Performance Cluster (HPC). The standard compute nodes were used for feature extraction and model training, each of which consists of 192 GB of RAM and 40 CPU cores at 2.50 GHz. The models (GMM, SVM, CNN, and CNN-GRU) were trained and tested utilizing the HPC's GPU nodes, made up of four NVIDIA Tesla V100 16G GPUs, 192 GB of RAM, and 48 CPU cores running at 2.10 GHz. The VLFEAT@matlab API was used for the GMMs, which consisted of 512 Gaussian clusters, and the GMMs were trained for 10 iterations.

## IX. Experimental Results

In this section, we report the results of the experiments performed to determine the effectiveness of the SOTA countermeasures. The effectiveness of each countermeasure was tested against standalone as well as cross-corpus evaluation. Standalone performance analysis implies the countermeasure's capability of detecting the spoofing artifacts extant in the training samples and similar ones existing in the testing samples. In contrast, cross-corpus evaluation determines the performance of each countermeasure against unknown and distinctively configured spoofing artifacts. For a fair comparison of the countermeasures, we test their effectiveness against four different ML and DL models (GMM, SVM, CNN, and CNN-GRU) and two diverse datasets, VSDC and ASVspoof2019. The VSDC dataset is used to evaluate the performance of countermeasures on smart home-based ASV systems, where microphonic discrepancies are more frequent, while the countermeasures are also tested against the ASVspoof19 industry standard spoofing discrepancies. Lastly, we comprehensively discuss the overall performance, both standalone and cross-corpus, of the countermeasures on all datasets and classifiers.

### A. Comparative Analysis of SOTA countermeasures on GMM based classifier

In this experiment, we evaluate the performance of the SOTA countermeasure with the baseline GMM classifier provided by the ASVspoof [103] challenge community. Although the classifier is initially provided with LFCC and CQCC features, we expanded it to evaluate the SOTA countermeasures. The experimental results show that each countermeasure performs significantly differently for speech synthesis (SS), VC, and replay attacks. The results demonstrate that the samples with replay attacks are much more difficult to detect in comparison with the SS and VC samples. During the evaluation, it is observed that the performance of the SOTA countermeasures declines drastically in the presence of replay artifacts in the speech samples.

The countermeasures were considered and analyzed for LA and PA attacks in the ASVspoof19 dataset examination. The countermeasures were tested against the development and evaluation speech samples, and the results are reported in Table XIII. The results demonstrated that each countermeasure achieved a good EER and min-TDCF against the development set of each spoofing attack (LA and PA). In particular, the SCMC [44] countermeasure achieved a lower EER of 0.01, while the lowest minimum t-DCF was obtained by the IM-FCC [9] and APGDF [67] countermeasures in the testing of development samples. However, in evaluation set testing, the APGDF [67] countermeasure achieved a lower EER and m-TDCF of 5.75 and 0.14, respectively. The results demonstrated that the APGDF [67] and CQCC [83] countermeasures outperformed the comparatives with an overall lower EER and min-TDCF in the development and evaluation testing XIII. In addition, some countermeasures performed optimally against the LA attacks comprised of SS and VC but failed to perform optimally when replay artifacts existed in the speech sample. For instance, the SCMC [44] countermeasures obtained the

**TABLE XIII** Experimental performance of the countermeasures with ASVspoof2019 and GMM based classifier.

| Feature | Dataset | Development | | Evaluation | |
|---------|---------|-------------|------|-------------|------|
| | | EER | m-tdcf | EER | m-tdcf |
| LFCC [113] | ASVspoof-LA | 2.70 | 0.06 | 8.08 | 0.21 |
| | ASVspoof-PA | 11.9 | 0.25 | 13.5 | 0.30 |
| CQCC [83] | ASVspoof-LA | 0.43 | 0.02 | 9.57 | 0.23 |
| | ASVspoof-PA | 4.87 | 0.19 | 11.0 | 0.24 |
| LPCC [91] | ASVspoof-LA | 2.31 | 0.06 | 10.17 | 0.28 |
| | ASVspoof-PA | 38.5 | 0.83 | 46.6 | 0.98 |
| MSRCC [80] | ASVspoof-LA | 9.45 | 0.19 | 10.9 | 0.28 |
| | ASVspoof-PA | 12.7 | 0.28 | 15.7 | 0.38 |
| PSRCC [80] | ASVspoof-LA | 9.13 | 0.14 | 10.7 | 0.21 |
| | ASVspoof-PA | 12.7 | 0.28 | 15.7 | 0.38 |
| SCFC [44] | ASVspoof-LA | 14.7 | 0.39 | 20.6 | 0.54 |
| | ASVspoof-PA | 17.6 | 0.35 | 21.6 | 0.47 |
| SCMC [44] | ASVspoof-LA | 0.01 | 0.54 | 5.91 | 0.15 |
| | ASVspoof-PA | 12.5 | 0.27 | 13.9 | 0.33 |
| MFCC [58] | ASVspoof-LA | 7.06 | 0.16 | 10.56 | 0.25 |
| | ASVspoof-PA | 11.5 | 0.23 | 13.7 | 0.32 |
| IMFCC [9] | ASVspoof-LA | 0.04 | 0.01 | 10.9 | 0.24 |
| | ASVspoof-PA | 12.8 | 0.29 | 13.7 | 0.32 |
| RFCC [73] | ASVspoof-LA | 2.71 | 0.07 | 8.06 | 0.22 |
| | ASVspoof-PA | 11.8 | 0.26 | 13.9 | 0.33 |
| RPS [73] | ASVspoof-LA | 9.28 | 0.19 | 11.9 | 0.29 |
| | ASVspoof-PA | 15.3 | 0.32 | 14.0 | 0.36 |
| SSFC [80] | ASVspoof-LA | 8.10 | 0.16 | 10.3 | 0.27 |
| | ASVspoof-PA | 13.8 | 0.29 | 13.9 | 0.31 |
| GTCC [85] | ASVspoof-LA | 9.25 | 0.17 | 10.8 | 0.24 |
| | ASVspoof-PA | 12.7 | 0.28 | 15.7 | 0.38 |
| APGDF [67] | ASVspoof-LA | 0.22 | 0.01 | 5.75 | 0.14 |
| | ASVspoof-PA | 8.66 | 0.17 | 10.6 | 0.25 |

best EER of 0.01 and the second best EER of 5.91 for the development and evaluation test sets, however, SCMC failed to provide optimal performance in the case of PA attack testing. Similarly, the EER of APGDF [67] countermeasure was higher than SCMC in the development set, but APGDF [67] achieved the overall best results in the rest of the testing. The detailed results for SOTA countermeasures are reported in Table XIII.

Along with the ASVspoof2019, we also examined the performance of the SOTA countermeasure with the VSDC dataset, and the results are shown in Table XIV. In the instance of the VSDC dataset, we computed the EER to assess the effectiveness of the countermeasures. The results showed that the LPCC [91] countermeasure exhibited the best and CQCC [83] obtained the second best EER of 5.86 and 20.0 in the VSDC dataset testing, respectively. Although the VSDC dataset contains 1PR and 2PR replay samples, this experiment only included 1PR testing because the goal was to distinguish between spoofed and genuine samples.It was observed from the results that GTCC [85] features proved to be an expensive countermeasure. The GTCC countermeasure [85] obtained the highest EER of 55.0. The results demonstrate the inadequacy of the countermeasures in testing the speech sample with the artifacts of distinct microphonic discrepancies. Except for the LFCC [113] countermeasure, none of the countermeasures achieved an EER lower than 40.0. The detailed countermeasure results of this experiment are presented in Table XIV.

From this experiment, we can infer that the countermeasures are totally dependent on the dataset characteristics. In the ASVspoof2019 dataset, for instance, the APGDF [67] features

**TABLE XIV** Experimental performance of the countermeasures with VSDC and GMM based classifier.

| Features | APGDF [67] | IMFCC [9] | GTCC [85] | CQCC [83] | LFCC [113] | RPS [73] | PSRCC [80] |
|----------|-----------|-----------|-----------|-----------|-----------|----------|-----------|
| EER | 40.29 | 50.00 | 56.00 | 20.0 | 36.33 | 55.21 | 55.11 |
| Features | LPCC [91] | MFCC [58] | RFCC [73] | SCFC [44] | SCMC [44] | SSFC [80] | MSRCC [80] |
| EER | 5.86 | 49.40 | 48.87 | 45.49 | 49.91 | 53.11 | 54.73 |

produced the best results. On the VSDC dataset, however, the APGDF [67] features failed to generate effective outcomes. Similarly, LPCC [91] did not generate effective results for the ASVspoof datasets, particularly in the ASVspoof-PA assessment set, but produced the best results in VSDC testing. Furthermore, we can see that the front end features that worked well against logical attacks did not perform well against physical attacks. The VSDC dataset also includes PA (replay) attacks that resulted in the performance degradation of each countermeasure. However, despite the dataset's wide variability, the APGDF [67], CQCC [83], and LPCC [91] features produce comparable results independent of the dataset's specific configuration.

### B. Performance analysis of SOTA countermeasures with an SVM classifier

In this experiment, we evaluate the effectiveness of SOTA countermeasures against the SVM classifier. The SVM classifier has been shown to be a superior optimal classifier to the GMM classifier [55]. An SVM with an RBF distribution was utilized in the experiment. The RBF distribution was selected after a thorough analysis of all feasible distributions. The scikit-learn [66] machine learning library was used to train and test the SVM classifier. In addition, the classification report [66] was used to compute the precision, recall, F1-score, and accuracy of the countermeasures. The results of this experiment on the VSDC dataset are mentioned in Table XV, and demonstrate that the CQCC [83] countermeasure proved to have the most robust front-end attributes in comparison with the SOTA countermeasures. CQCC [83] features, in particular, achieved the lowest EER of 0.27, the highest precision, recall, and F1-scores of 0.92, 0.85, and 0.88, respectively, as well as the highest accuracy of 0.93. Spectral features such as SCFC [44], on the other hand, performed poorly in the case of the VSDC dataset, which contains audio samples with microphonic inconsistencies. For instance, the SCFC [44] features had the lowest precision (0.39), recall (0.50), and F1-score (0.44), and achieved an accuracy of 0.78. Similarly, the SCMC [44] features were found to be the second worst classifier, after the SCFC [44], with a 0.53 EER. This indicates that the spectral centroid magnitude and frequency-based coefficients failed to perform effectively when micro-phonic distinction (i.e., VSDC) exists in the speech sample. In contrast, the constant Q-cepstral-based coefficient successfully overcame the microphone's fluctuation.

In the case of the ASVspoof2019 dataset, each of the SOTA countermeasures was tested against the SVM classifier, along with the development and evaluation set of the datasets, and the results are reported in Table XVI. The results demonstrate that in the testing of LA attacks that exist in the development set and evaluation set, the CQCC [83] and IMFCC [9] coun-

termeasures achieved comparable performance, with an EER of 0.30 and 0.29, respectively. Similarly, for the evaluation set, the CQCC [83] and IMFCC [9] countermeasures achieved the best EER of 0.69 and 0.76, respectively. Although the EER of 0.69 and 0.76 was not good enough for any spoofing system, none of the other countermeasures achieved an EER lower than these countermeasures. From this experiment, it was concluded that the performance of the majority of the spectral and cepstral countermeasures declined when tested against different classifiers with different configurations. However, the CQCC [83] countermeasure performed far better than the other countermeasures. In addition, in the presence of micro-phonic distinctions, CQCC [83] features work effectively, while the rest of the spectral features fail to perform adequately. Furthermore, the phase-based APGDF [67] features were the third best features compared to the state-of-the-art that achieved the optimal results. Detailed results of the SOTA countermeasures are presented in Table XVI.

Based on the experiments involving the SOTA countermeasures with machine learning-based GMM and SVM-based classifiers, it is determined that the CQCC [83] is the best countermeasure against voice spoofing attacks. More specifically, in the testing of replay attacks, where the performance of the state-of-the-art classifiers drops significantly, the CQCC [83] features scored the optimal EER of 0.27, 0.43, and 0.30 in GMM and SVM testing, respectively. The CQCC [83] countermeasure achieved the lowest EER and min-TDCF when in development. The results of both ML-based classifiers show that the phase-based APGDF [67], linear LFCC [113] , and inverted Mel-based IMFCC [9] countermeasures performed better than the comparative methods in all experiments and received the optimal scores in all performance measurements. APGDF [67] also had the lowest EER and min-TDCF in the ASVspoof2019-PA dataset evaluation. Similar results are obtained with the VSDC and SVM classifiers. In comparison to the comparative approaches, the SCFC [44] features were shown to be the weakest. SCFC [44] failed to perform well in any experiment, and had the highest EER and min-TDCF. Although the CQCC [83] and APGDF [67] front end features clearly outperformed all comparable features in the ASVspoof2019-LA, ASVspoof2019-PA, and VSDC datasets, the question of performance when the ASV encounters unfamiliar samples from unknown surroundings and devices persists. Although the ML-based GMM and SVM classifiers showed the capabilities of the countermeasure, they have become the conventional classifiers for spoof detection. In contrast, several deep learning and end-to-end solutions have been published recently to differentiate between spoofed and bonafide speech samples. Therefore, we tested the countermeasures against the recent CNN and CNN-GRU based classifiers and also performed cross-corpus experiments in

**TABLE XV** Experimental performance of the countermeasures with VSDC dataset and SVM based classifier.

| Feature | EER | Precision | Recall | F1-score | Accuracy |
|---|---|---|---|---|---|
| LFCC [113] | 0.49 | 0.88 | 0.74 | 0.79 | 0.88 |
| CQCC [83] | 0.27 | 0.92 | 0.85 | 0.88 | 0.93 |
| LPCC [91] | 0.44 | 0.88 | 0.77 | 0.81 | 0.89 |
| PSRCC [80] | 0.40 | 0.79 | 0.76 | 0.77 | 0.86 |
| MSRCC [80] | 0.46 | 0.84 | 0.76 | 0.79 | 0.86 |
| SCFC [44] | 0.80 | 0.39 | 0.50 | 0.44 | 0.78 |
| SCMC [44] | 0.53 | 0.87 | 0.72 | 0.76 | 0.87 |
| MFCC [58] | 0.46 | 0.82 | 0.75 | 0.77 | 0.86 |
| IMFCC [9] | 0.71 | 0.88 | 0.64 | 0.67 | 0.84 |
| RPS [73] | 0.40 | 0.76 | 0.78 | 0.77 | 0.85 |
| RFCC [73] | 0.51 | 0.87 | 0.73 | 0.76 | 0.87 |
| GTCC [85] | 0.42 | 0.82 | 0.76 | 0.78 | 0.86 |
| APGDF [67] | 0.43 | 0.89 | 0.77 | 0.81 | 0.89 |
| SSFC [80] | 0.44 | 0.87 | 0.77 | 0.80 | 0.88 |

**TABLE XVI** Experimental performance of the countermeasures with ASVspoof2019 and an SVM based classifier.

| Feature | Dataset | Development | | Evaluation | |
|---|---|---|---|---|---|
| | | EER | Acc | EER | Acc |
| LFCC [113] | ASVspoof19-LA | 0.58 | 0.90 | 0.91 | 0.12 |
| | ASVspoof19-PA | 0.72 | 0.75 | 0.83 | 0.11 |
| CQCC [83] | ASVspoof19-LA | 0.30 | 0.94 | 0.69 | 0.29 |
| | ASVspoof19-PA | 0.42 | 0.29 | 0.80 | 0.19 |
| LPCC [91] | ASVspoof19-LA | 0.50 | 0.88 | 0.80 | 0.15 |
| | ASVspoof19-PA | 0.61 | 0.20 | 0.81 | 0.19 |
| SCFC [44] | ASVspoof19-LA | 0.50 | 0.89 | 0.90 | 0.16 |
| | ASVspoof19-PA | 0.75 | 0.19 | 0.89 | 0.14 |
| SCMC [44] | ASVspoof19-LA | 0.49 | 0.93 | 0.88 | 0.12 |
| | ASVspoof19-PA | 0.75 | 0.20 | 0.86 | 0.15 |
| PSRCC [80] | ASVspoof19-LA | 0.49 | 0.89 | 0.89 | 0.11 |
| | ASVspoof19-PA | 0.84 | 0.25 | 0.89 | 0.12 |
| MSRCC [80] | ASVspoof19-LA | 0.51 | 0.90 | 0.91 | 0.11 |
| | ASVspoof19-PA | 0.84 | 0.25 | 0.89 | 0.12 |
| MFCC [58] | ASVspoof19-LA | 0.50 | 0.89 | 0.80 | 0.17 |
| | ASVspoof19-PA | 0.76 | 0.22 | 0.86 | 0.16 |
| IMFCC [9] | ASVspoof19-LA | 0.29 | 0.94 | 0.76 | 0.24 |
| | ASVspoof19-PA | 0.51 | 0.26 | 0.82 | 0.18 |
| RFCC [73] | ASVspoof19-LA | 0.46 | 0.93 | 0.87 | 0.14 |
| | ASVspoof19-PA | 0.84 | 0.26 | 0.88 | 0.14 |
| RPS [73] | ASVspoof19-LA | 0.52 | 0.90 | 0.91 | 0.10 |
| | ASVspoof19-PA | 0.86 | 0.25 | 0.89 | 0.11 |
| SSFC [80] | ASVspoof19-LA | 0.50 | 0.90 | 0.89 | 0.12 |
| | ASVspoof19-PA | 0.87 | 0.29 | 0.90 | 0.10 |
| GTCC [85] | ASVspoof19-LA | 0.50 | 0.90 | 0.90 | 0.10 |
| | ASVspoof19-PA | 0.84 | 0.25 | 0.89 | 0.12 |
| APGDF [67] | ASVspoof19-LA | 0.34 | 0.94 | 0.82 | 0.17 |
| | ASVspoof19-PA | 0.60 | 0.17 | 0.83 | 0.15 |

order to validate the generalized performance of both the highest performing countermeasures. The results are discussed in the next subsection.

### C. Performance analysis of the SOTA countermeasures with a CNN-based classifier

In this subsection, we show the effectiveness of SOTA countermeasures along with the DL-based CNN and CNN-GRU classifier. The CNN classifier is one of the most recent and most advanced neural network-based classifiers for the detection of voice spoofing on ASV systems. The SOTA countermeasures were tested against the CNN classifier, and the results are shown in Table XVII. The same architectural CNN

was used in the experiments that was adopted in [34]. The open source code is implemented in the mentioned configuration in section VIII subsection C.

The results for the VSDC dataset demonstrate that the APGDF [67] front end features outperform the SOTA countermeasures with an EER of 0.01 and a min T-DCF of 0.04. Following all of the countermeasures, the linear-based features perform well with the CNN-based classifier to detect voice spoofing attacks. For instance, according to performance analysis with EER, LPCC [91] and LFCC [113] rank second and third, with EERs of 0.026 and 0.032, respectively. Similarly, in the case of ASVspoof19, the rectangular-based RFCC [73] features outperform all evaluated countermeasures on the ASVspoof19 dataset, with an EER of 0.054 and a min T-DCF of 0.139. The magnitude-based feature was found to be significantly deficient during the ASVspoof19 evaluation. For instance, MSRCC [80] features had the highest EER and min-TDCF values, at 0.443 and 0.947, respectively. In contrast, the highest EER of 0.21, and the highest min-TDCF score of 0.43, were seen in the SCFC [44] features against VSDC. Moreover, when compared to the GMM classifier, the CNN classifier performed significantly better across the board. For instance, the highest EER recorded with the GMM classifier was 0.80 on the LFCC feature, while the highest EER achieved with the CNN classifier was just 0.265 with the MFCC [58] features. Using the same feature set but replacing the classifier, the EER was reduced by nearly 65%. This experiment clearly shows the significance of back end classifiers in speech spoofing detection. Detailed results of the SOTA countermeasures are presented in Table XVII.

**TABLE XVII** Experimental performance analysis of the countermeasures on the training and testing subsets of the same dataset with a CNN based classifier.

| Features | Training | Testing | EER | min-tdcf |
|---|---|---|---|---|
| IMFCC [9] | VSDC | VSDC | 0.06 | 0.17 |
| | ASVspoof19 | ASVspoof19 | 0.19 | 0.50 |
| SSFC[80] | VSDC | VSDC | 0.03 | 0.07 |
| | ASVspoof19 | ASVspoof19 | 0.18 | 0.46 |
| RFCC [73] | VSDC | VSDC | 0.13 | 0.25 |
| | ASVspoof19 | ASVspoof19 | 0.05 | 0.13 |
| SCFC [44] | VSDC | VSDC | 0.10 | 0.30 |
| | ASVspoof19 | ASVspoof19 | 0.14 | 0.36 |
| SCMC[44] | VSDC | VSDC | 0.14 | 0.37 |
| | ASVspoof19 | ASVspoof19 | 0.44 | 0.94 |
| PSRCC [80] | VSDC | VSDC | 0.08 | 0.10 |
| | ASVspoof19 | ASVspoof19 | 0.43 | 0.75 |
| MSRCC [80] | VSDC | VSDC | 0.13 | 0.17 |
| | ASVspoof19 | ASVspoof19 | 0.43 | 0.75 |
| RPS [73] | VSDC | VSDC | 0.04 | 0.11 |
| | ASVspoof19 | ASVspoof19 | 0.21 | 0.65 |
| CQCC [83] | VSDC | VSDC | 0.04 | 0.13 |
| | ASVspoof19 | ASVspoof19 | 0.17 | 0.45 |
| LFCC [113] | VSDC | VSDC | 0.03 | 0.03 |
| | ASVspoof19 | ASVspoof19 | 0.16 | 0.43 |
| MFCC [58] | VSDC | VSDC | 0.06 | 0.16 |
| | ASVspoof19 | ASVspoof19 | 0.26 | 0.62 |
| LPCC [91] | VSDC | VSDC | 0.02 | 0.05 |
| | ASVspoof19 | ASVspoof19 | 0.19 | 0.50 |
| GTCC [85] | VSDC | VSDC | 0.03 | 0.07 |
| | ASVspoof19 | ASVspoof19 | 0.43 | 0.75 |
| APGDF [67] | VSDC | VSDC | 0.01 | 0.04 |
| | ASVspoof19 | ASVspoof19 | 0.17 | 0.45 |

**TABLE XVIII** Experimental performance analysis of the SOTA countermeasures with a CNN-GRU based classifier.

| Feature | Train Dataset | Test Dataset | EER | min T-DCF |
|---|---|---|---|---|
| LFCC [113] | VSDC | VSDC | 0.02 | 0.05 |
| | ASVspoof19 | ASVspoof19 | 0.15 | 0.41 |
| SSFC[80] | VSDC | VSDC | 0.04 | 0.07 |
| | ASVspoof19 | ASVspoof19 | 0.17 | 0.46 |
| MFCC [58] | VSDC | VSDC | 0.03 | 0.07 |
| | ASVspoof19 | ASVspoof19 | 0.22 | 0.548 |
| SCFC[44] | VSDC | VSDC | 0.18 | 0.35 |
| | ASVspoof19 | ASVspoof19 | 0.09 | 0.25 |
| LPCC [91] | VSDC | VSDC | 0.18 | 0.03 |
| | ASVspoof19 | ASVspoof19 | 0.02 | 0.47 |
| SCMC[44] | VSDC | VSDC | 0.02 | 0.05 |
| | ASVspoof19 | ASVspoof19 | 0.17 | 0.46 |
| MSRCC [80] | VSDC | VSDC | 0.21 | 0.03 |
| | ASVspoof19 | ASVspoof19 | 0.22 | 0.50 |
| PSRCC [80] | VSDC | VSDC | 0.23 | 0.03 |
| | ASVspoof19 | ASVspoof19 | 0.19 | 0.53 |
| RFCC [73] | VSDC | VSDC | 0.11 | 0.22 |
| | ASVspoof19 | ASVspoof19 | 0.04 | 0.11 |
| IMFCC[9] | VSDC | VSDC | 0.03 | 0.09 |
| | ASVspoof19 | ASVspoof19 | 0.18 | 0.50 |
| APGDF [67] | VSDC | VSDC | 0.08 | 0.26 |
| | ASVspoof19 | ASVspoof19 | 0.14 | 0.37 |
| CQCC [83] | VSDC | VSDC | 0.09 | 0.23 |
| | ASVspoof19 | ASVspoof19 | 0.18 | 0.49 |
| GTCC [85] | VSDC | VSDC | 0.04 | 0.09 |
| | ASVspoof19 | ASVspoof19 | 0.19 | 0.53 |
| RPS [73] | VSDC | VSDC | 0.04 | 0.09 |
| | ASVspoof19 | ASVspoof19 | 0.19 | 0.53 |

### D. Experimental Analysis of the SOTA countermeasures with a CNN-GRU based classifier

In this experiment, a CNN-GRU [34] classifier was used to classify the SOTA countermeasures in order to test their effectiveness against the SS, VC, and replay attacks, and the results of this experiment are presented in Table XVIII. The results demonstrate that in the case of the VSDC dataset, the linear and spectral magnitude-based LFCC and SCMC features achieved a comparatively similar performance, with an EER of 0.02 and a min T-DCF of 0.05. Due to the versatility of the VSDC dataset, the SCFC features scored an EER of 0.18 and a min T-DCF of 0.35. In particular, the results from this experiment proved to be optimal for the VSDC dataset, in comparison with the GMM, SVM, and CNN-GRU classifiers. Moreover, the SOTA countermeasures achieved the best result on the CNN-GRU classifiers. In contrast, in the case of ASVspoof19, LPCC was the best performing countermeasure, with an EER of 0.02 and a min T-DCF of 0.47, while the MFCC and MSRCC countermeasures have the lowest EERs of 0.22 and 0.18, respectively. Therefore, it was concluded from this experiment that the performance of the existing SOTA countermeasure differs in the presence of distinct classifiers and may not function adequately in the presence of advanced spoofing attempts. In addition, we can infer that the CNN-GRU classifier proved to be the optimal classifier for both of the datasets. However, the effectiveness of the SOTA countermeasure across corpora must still be demonstrated. Therefore, in the next section we provide a cross-corpus evaluation of the SOTA countermeasures.

### E. Cross-corpus evaluations of SOTA countermeasures

The SOTA countermeasures are compared in this section using cross-corpus evaluation. Cross-corpus evaluation and generalisation was one of the key needs mentioned by the in-field researcher [79], which we resolved in this experiment. A key shortcoming in previous research is the lack of cross-corpus experiments [69]. Cross-corpus examination demonstrates the countermeasure's effectiveness in dealing with unknown aspects of spoofing attacks by removing prior bias from the model in the form of characteristics shared by a particular dataset, and allows it to perform a classification based solely on how the model recognises the unique aspects of spoofing attacks. In addition, this experiment was performed to determine the comprehensiveness of the countermeasures to speech samples, including unknown and diverse spoofing artifacts such as distinct VC algorithms, SSD techniques, and others. During the experiment, the countermeasures were trained using one dataset and the testing was done using the speech samples from another. There are some similarities between the two datasets. Replayed speech samples, for instance, can be found in the VSDC and ASVspoof19 databases. Initially, the countermeasures were trained using VSDC data samples and then tested with the replayed spoof speech samples from the ASVspoof corpus. The SOTA countermeasure was then trained using ASVspoof and evaluated using VSDC data samples. This evaluation was performed against three distinct classifiers, and the results are shown below in Tables XIX, XX, and XXI.

#### 1) Cross-corpus evaluation of the SOTA countermeasures using a GMM classifier

In this subsection, performance of the SOTA countermeasure was against a GMM classifier, and the results are reported in Table XIX. The results demonstrate that when the classifier was trained on the VSDC dataset and tested on the ASVspoof19 dataset, the countermeasure PSRCC [80] outperformed the SOTA countermeasures, with an EER of 0.25 and a min TDCF of 0.39. When trained on the ASVspoof19 dataset and evaluated on the VSDC dataset, RFCC [73] was the top performing countermeasure, with the lowest EER of 0.39 and the lowest min T-DCF of 0.27. In contrast, when training on VSDC and testing with the ASVspoof dataset, the countermeasure SSFC [44] obtained the highest EER of 0.80 and min-tdcf of 0.99. Similarly, the SSFC [80] countermeasure obtained the highest EER of 0.67. It is observed from the results that the EER of the SOTA countermeasures varies from 0.25 to 0.80, but min-TDCF ranges from 0.81 to 1.0. Only the PSRCC [80], MSRCC[80], RFCC [73], GTCC [85], and APGDF [67] countermeasures are able to achieve lower min-TDCF, at 0.39, 0.47, 0.27, 0.38, and 0.59, respectively. The EER performance metric is used to evaluate the effectiveness of spoofing countermeasures, whereas the min-TDCF evaluation parameter can be used to evaluate the efficacy of ASVs. Consequently, based on these results, we may conclude that the effectiveness of the SOTA countermeasures will be called into question when employed for ASV evaluations.

Surprisingly, no specific SOTA countermeasure performed considerably better in both cross-corpus evaluations. For instance, when trained using ASVspoof19 and evaluated against

the VSDC dataset, the countermeasure MSRCC worked effectively. In the ASVspoof dataset test, EER increased by up to 11%. Despite being the worst performer in the standalone tests, PSRCC performed well when trained on VSDC and evaluated with ASVspoof19. While phase-based spectral features outperformed the state-of-the-art in cross-corpus evaluation, the current investigation revealed that when utilizing the GMM classifier, the performance of the existing countermeasures differed greatly in cross-corpus evaluation. Moreover, the results show that when the system was unfamiliar with the testing speech samples, the countermeasures failed to accurately recognize the spoofed speech samples. However, determining the countermeasure's cross-corpus performance using a single classifier is not practical for fairness and may be biased. Therefore, along with the GMM classifier, we demonstrate the efficacy of the SOTA countermeasures against the DL-based CNN and CNN-GRU classifiers, and the results are reported below.

*2) Cross corpus evaluation of SOTA countermeasures with a CNN based classifier*

In this section, cross-corpus evaluations of the SOTA countermeasures with a CNN-based classifier are discussed, and the results are presented in Table XX. The results demonstrate that in VSDC training and ASVspoof19 testing, the phase-oriented PSRCC [80] countermeasure performed the best and obtained an EER of 0.37 and a min T-DCF of 0.94. In contrast, in the case of ASVspoof19 training and VSDC testing, the SCMC [44] features achieved the minimum values for EER and min T-DCF, at 0.35 and 0.96, respectively. Although the countermeasures and the CNN-based classifier performed marginally better when trained with VSDC and tested with

**TABLE XIX** Cross corpus evaluation of the SOTA countermeasures with GMM.

| Features | Training | Testing | EER | min-tdcf |
|---|---|---|---|---|
| CQCC [83] | VSDC | ASVspoof19 | 0.40 | 0.81 |
| | ASVspoof19 | VSDC | 0.66 | 0.99 |
| LFCC [113] | VSDC | ASV'19 PA | 0.60 | 0.99 |
| | ASV'19 PA | VSDC | 0.45 | 0.99 |
| SSFC [80] | VSDC | ASVspoof19 | 0.80 | 0.99 |
| | ASVspoof19 | VSDC | 0.67 | 0.99 |
| RPS [73] | VSDC | ASVspoof19 | 0.50 | 0.93 |
| | ASVspoof19 | VSDC | 0.40 | 0.94 |
| LPCC [91] | VSDC | ASVspoof19 | 0.41 | 0.89 |
| | ASVspoof19 | VSDC | 0.62 | 1.0 |
| PSRCC [80] | VSDC | ASVspoof19 | 0.25 | 0.39 |
| | ASVspoof19 | VSDC | 0.67 | 0.99 |
| MSRCC [80] | VSDC | ASVspoof19 | 0.29 | 0.47 |
| | ASVspoof19 | VSDC | 0.40 | 0.94 |
| SCFC [44] | VSDC | ASVspoof19 | 0.58 | 0.99 |
| | ASVspoof19 | VSDC | 0.48 | 0.99 |
| SCMC[44] | VSDC | ASVspoof19 | 0.60 | 0.97 |
| | ASVspoof19 | VSDC | 0.46 | 0.87 |
| MFCC [58] | VSDC | ASVspoof19 | 0.42 | 0.97 |
| | ASVspoof19 | VSDC | 0.45 | 0.93 |
| IMFCC [9] | VSDC | ASVspoof19 | 0.45 | 0.99 |
| | ASVspoof19 | VSDC | 0.45 | 0.96 |
| RFCC [73] | VSDC | ASVspoof19 | 0.34 | 0.84 |
| | ASVspoof19 | VSDC | 0.39 | 0.27 |
| GTCC [85] | VSDC | ASVspoof19 | 0.36 | 0.38 |
| | ASVspoof19 | VSDC | 0.56 | 1.0 |
| APGDF [67] | VSDC | ASVspoof19 | 0.30 | 0.59 |
| | ASVspoof19 | VSDC | 0.61 | 0.99 |

**TABLE XX** Cross corpus evaluation of the SOTA countermeasures with a CNN based classifier.

| Feature | Train Dataset | Test Dataset | EER | min T-DCF |
|---|---|---|---|---|
| IMFCC [9] | VSDC | ASVspoof19 | 0.43 | 0.99 |
| | ASVspoof19 | VSDC | 0.57 | 1.0 |
| SCMC [44] | VSDC | ASVspoof19 | 0.51 | 0.99 |
| | ASVspoof19 | VSDC | 0.35 | 0.96 |
| RFCC [73] | VSDC | ASVspoof19 | 0.40 | 0.96 |
| | ASVspoof19 | VSDC | 0.50 | 0.91 |
| RPS [73] | VSDC | ASVspoof19 | 0.68 | 0.99 |
| | ASVspoof19 | VSDC | 0.55 | 0.98 |
| SSFC [80] | VSDC | ASVspoof19 | 0.72 | 0.99 |
| | ASVspoof19 | VSDC | 0.39 | 0.91 |
| PSRCC [80] | VSDC | ASVspoof19 | 0.37 | 0.94 |
| | ASVspoof19 | VSDC | 0.51 | 0.97 |
| MSRCC [80] | VSDC | ASVspoof19 | 0.46 | 0.52 |
| | ASVspoof19 | VSDC | 0.47 | 0.93 |
| CQCC [83] | VSDC | ASVspoof19 | 0.48 | 0.98 |
| | ASVspoof19 | VSDC | 0.66 | 1.0 |
| LFCC [113] | VSDC | ASVspoof19 | 0.51 | 0.99 |
| | ASVspoof19 | VSDC | 0.37 | 0.99 |
| MFCC [58] | VSDC | ASVspoof19 | 0.47 | 0.93 |
| | ASVspoof19 | VSDC | 0.52 | 0.94 |
| SCFC[44] | VSDC | ASVspoof19 | 0.53 | 0.99 |
| | ASVspoof19 | VSDC | 0.54 | 0.97 |
| LPCC [91] | VSDC | ASVspoof19 | 0.46 | 0.98 |
| | ASVspoof19 | VSDC | 0.63 | 1.0 |
| GTCC [85] | VSDC | ASVspoof19 | 0.50 | 0.98 |
| | ASVspoof19 | VSDC | 0.41 | 0.91 |
| APGDF [67] | VSDC | ASVspoof19 | 0.44 | 0.99 |
| | ASVspoof19 | VSDC | 0.56 | 1.0 |

the ASVspoof19 dataset, several of the front-end features performed significantly worse. This uncertainty shows that the performance of these countermeasures remains ambiguous against unseen attacks without having prior knowledge of the speech sample.

*3) Cross corpus measurements of the SOTA countermeasures with a CNN-GRU based classifier*

A CNN-GRU classifier was used to evaluate the countermeasures across corpora in this subsection, and the results are shown in Table XX. In the instance of VSDC training and ASVspoof19 testing, the APGDF-based countermeasure [67] obtained an optimal EER of 0.34 and a minimum T-DCF score of 0.82. With an EER of 0.44, the RFCC [73] and IMFCC [9] were ranked second and third, respectively. In the case of ASVspoof19 training and VSDC testing, the MFCC-based countermeasure [58] outperformed the SOTA countermeasures with the lowest EER of 0.34 and the lowest min T-DCF score of 0.97. The outcomes of this experiment were similar to the results of the CNN-based cross-corpus evaluations. The min-TDCF of the SOTA countermeasures was greater than the countermeasures' standalone testing. In the instance of ASVspoof19 dataset training and VSDC testing, the results demonstrated that the majority of the countermeasures' minimum T-DCF was at or around 1. These findings demonstrate the ineffectiveness of existing countermeasures and classifiers against unknown spoofing attacks. Furthermore, these findings have raised concerns regarding the effectiveness of existing protections against complicated contemporary spoofing attacks.

**TABLE XXI** Cross corpus evaluation of the countermeasures with a CNN-GRU based classifier.

| Feature | Train Dataset | Test Dataset | EER | min T-DCF |
|---|---|---|---|---|
| LFCC [113] | VSDC | ASVspoof19 | 0.47 | 0.99 |
| | ASVspoof19 | VSDC | 0.53 | 0.98 |
| MFCC [58] | VSDC | ASVspoof19 | 0.45 | 0.94 |
| | ASVspoof19 | VSDC | 0.34 | 0.97 |
| SCFC[44] | VSDC | ASVspoof19 | 0.49 | 0.99 |
| | ASVspoof19 | VSDC | 0.59 | 0.99 |
| LPCC [91] | VSDC | ASVspoof19 | 0.65 | 0.99 |
| | ASVspoof19 | VSDC | 0.45 | 1.0 |
| SCMC[44] | VSDC | ASVspoof19 | 0.45 | 0.99 |
| | ASVspoof19 | VSDC | 0.65 | 1.0 |
| PSRCC [80] | VSDC | ASVspoof19 | 0.521 | 0.987 |
| | ASVspoof19 | VSDC | 0.57 | 0.99 |
| MSRCC [80] | VSDC | ASVspoof19 | 0.49 | 0.99 |
| | ASVspoof19 | VSDC | 0.54 | 1.0 |
| RFCC [73] | VSDC | ASVspoof19 | 0.44 | 0.99 |
| | ASVspoof19 | VSDC | 0.67 | 1.0 |
| IMFCC [9] | VSDC | ASVspoof19 | 0.44 | 0.99 |
| | ASVspoof19 | VSDC | 0.61 | 1.0 |
| APGDF [67] | VSDC | ASVspoof19 | 0.34 | 0.82 |
| | ASVspoof19 | VSDC | 0.51 | 0.96 |
| CQCC [83] | VSDC | ASVspoof19 | 0.50 | 0.99 |
| | ASVspoof19 | VSDC | 0.43 | 1.0 |
| GTCC [85] | VSDC | ASVspoof19 | 0.49 | 0.99 |
| | ASVspoof19 | VSDC | 0.54 | 1.0 |
| SSFC [80] | VSDC | ASVspoof19 | 0.51 | 0.99 |
| | ASVspoof19 | VSDC | 0.56 | 1.0 |
| RPS [73] | VSDC | ASVspoof19 | 0.50 | 0.96 |
| | ASVspoof19 | VSDC | 0.55 | 0.98 |

## F. Overall performance of the SOTA countermeasures and a way forward

One of the shortfalls in the current research into audio spoof detection is that cross-corpus evaluations are rarely done, limiting knowledge of the generalizability of the features and full countermeasure frameworks. To overcome this knowledge gap, this comparative analysis performs a cross-corpus evaluation, where one dataset is used for training and another for evaluation, to show the effectiveness of these features at generalized spoof detection across different conditions present in the respective datasets.

Based on the results of this study's experiments, we can infer that the no-countermeasures from the selected ones were consistent over a wide range of datasets and classifiers. As shown in the preceding section, several countermeasures, such as CQCC [83] and RFCC [73], performed better overall in standalone testing of the countermeasures. However, when the dataset and classifier were changed, the identical countermeasure failed to perform well. While PSRCC [80] and MSRCC [80] did not perform much better in standalone testing, they did significantly better in cross-corpus examination. After extensive testing, it is extremely difficult to select the best countermeasure from among those tested because some countermeasures performed well with one data set but failed with the other, while others performed in reverse with other classifiers. However, we may infer that classifiers and datasets have a significant influence on the present countermeasures. In addition, in standalone corpus testing with four alternative classifiers, the constant Q based CQCC [83] performed consistently with only a minor decline in performance across classifiers. The countermeasures PSRCC [80] and APGDF [67]

scored marginally better than each other when the classifier varied and in cross-corpus evaluation. As a consequence of the extensive testing with various datasets and classifiers, the effectiveness of the countermeasures was called into question. It also demonstrates the significance of cross-corpus and variable classifier evaluation of proposed solutions.

## G. Performance Evaluation of recent deep learned and End-to-End solutions

Following the observation of a performance drop in SOTA countermeasures, we conducted a cross-corpus sub-experiment utilizing contemporary deep learning, along with an end-to-end solution. In this subsection, the performance of recent deep learning and end-to-end solutions is evaluated. The results are in Table XXII.

Several deep learning and end-to-end solutions have recently demonstrated outstanding performance in the recent ASV challenges [33], [78], [110], [2], [49]. In the latest ASVspoof2021 competition, for instance, the end-to-end solution ASSIST [33] achieved the best EER and min-TDCF on the ASVspoof2021-LA dataset. Similarly, [110] and [78] also performed well in this task. Many participants, however, presented deep learning and end-to-end systems that had been trained and tested to mitigate only a single attack. To demonstrate the general effectiveness of the countermeasures, we performed an experiment in which we evaluated the higher performing solutions to the ASV challenges across corpora. Each solution trained and tested on ASVspoof2019-LA for the challenge, evaluated against the LA speech samples of the ASVspoof2021 corpus. Similarly, if the model's stated results were based on ASVspoof2021, we examined that model's performance against ASVspoof2019. The key reasoning behind this experiment was to evaluate and report the generalizability of the countermeasures.

To measure generalizability, we selected the top-performing models with publicly available code. It was observed from the results that the top performer, ASSIST [33], and its companion model ASSIST-L, achieved EER values of 0.83% and 0.99%, respectively, against the ASVspoof2019 dataset. However, the performance of these models deteriorated when tested against speech samples from the ASVspoof2021 dataset. Specifically, the EER scores increased to 15.84% and 15.69%, and the minimum t-DCFs rose from 0.027% to 0.537%, and from 0.030% to 0.520%, respectively.

When tested against ASV spoof 2021, the EER for the solution proposed in [110] increased from 2.17% to 23.1%, and the minimum TDCF increases from 0.055% to 0.759%, however the worst performance was observed in the case of the RawBoost solution proposed in [78]. Although the proposed solution performed optimally with the augmented data approach with EER score of 2.17%, the EER and minimum t-DCF increased dramatically when tested on other datasets and without augmentation. In the case of ASVspoof2021 and without data augmentation, EER increased by 43%, and the minimum t-DCF increased from 0.301% to 0.903%.

Although the adaptation of the Resnet-34 architecture published in [2] reported an EER of 5.32% against

**TABLE XXII** Computational details of the state-of-the-art countermeasures.

| Model | ASVspoof2019 | | ASVspoof2021 | | VSDC | |
|---|---|---|---|---|---|---|
| | EER | min t-DCF | EER | min t-DCF | EER | min t-DCF |
| ASSIST [33] | 0.83 | 0.027 | 15.84 | 0.537 | 51.21 | – |
| ASSIST-L [33] | 0.99 | 0.030 | 15.69 | 0.520 | 59.32 | – |
| RawBoost [78] | 2.17 | 0.055 | 23.1 | 0.759 | 64.90 | – |
| One class Learning [110] | 5.38 | 0.309 | 48.7 | 0.901 | 37.95 | – |
| Res2Net34 [49] | 1.67 | 0.058 | 48.00 | 0.994 | 29.00 | – |
| Res2Net50 [49] | 1.63 | 0.049 | 39.00 | 0.999 | 24.00 | – |
| ResNet34-LA [2] | 10.53 | 0.19 | 26.55 | 0.34 | 40.23 | – |
| ResNet34-PA [2] | 6.61 | 0.15 | 49.51 | 0.41 | 34.12 | – |

the ASVSpoof2019-LA dataset and 5.74% against the ASVSpoof2019-PA dataset, we observed EERs of 10.53% and 6.61%, respectively, in our configured environment. When tested on ASVSpoof2021, performance dropped even further, with EER results of 26.55% on LA and 49.51% on PA.

To further test the extensibility of this and other methods, we performed cross-attack, cross-corpus testing, training on ASVspoof2019-PA, and testing against ASVSpoof2021-DF, which achieved an EER 42.21%. Similar performance deterioration was observed for [2] on the distinct dataset; the EER rose to 48.00% and 39.00% for res2net34 and res2net50, respectively. The performance decreased when tested on the ASVspoof2021 dataset, and the worst performance was observed against the chained replay VSDC datasets. These scores imply that while these methods perform well on a specific dataset, the results do not transfer well to others. This test demonstrates the impact of formulating a solution based on a particular corpus, as the performance of deep learning and end-to-end solutions degrades significantly when presented with unknowns. Each static dataset differs and has distinct discrepancies due to varied factors, such as the ambient setting, gadgets, microphonic influence, and the spoofing technique used to produce spoof samples. Therefore, developing a spoofing countermeasure without taking into consideration the aforementioned factors may lead to performance degradation when tested across corpora. Thus, based on these results, we conclude that there is a dire need to concentrate on a broad approach that mitigates the aforementioned variances while developing robust voice spoofing defenses.

### H. Limitation and Future work

Aside from the effectiveness of the presented work, there are certain limitations to this study. Due to contradictory technical issues and lack of code availability, only fourteen countermeasures were chosen to do the cross-corpus and multi-classifier evaluation. Moreover, in accordance with the advancement of the voice spoofing domain and deep learning methodologies such as CNNs, transformers, and so on, we compared four recent top-rated (by EER and min-TDCF scores on ASVspoof datasets) deep learning and end-to-end solutions with cross corpus evaluation. Despite the fact that just four models were examined, the findings motivate us to continue working in the deep learning and end-to-end arenas. Thus, we want to expand this work in the future to include the most recent deep learning/end-to-end solutions in order to provide an apples-to-apples comparison. Lastly, we intend to create a

benchmark for comparing CNN and DNN-based architectures for the automatic speaker verification system.

## X. Existing Issues and Future Needs

This section describes the challenges and issues observed during this review as well as the future requirements for developing dependable and secure anti-spoofing solutions.

### A. Need for Explainable AI in Audio Forensics

Explainable AI (XAI) is an emerging field which emphasises the need for human understanding through what has traditionally been a black box process. An important gain over the interaction between humans and XAI is a level of trust in the process, which studies have shown does not currently exist [20]. Existing deepfake detection approaches are typically designed to perform batch analysis over a large dataset, however, when these techniques are employed in the field by journalists or law enforcement, there may only be a small set of videos available for analysis. A numerical score parallel to the probability of an audio or video being real or fake is not valuable to these users if it cannot be confirmed with an appropriate proof of the score. In those situations, it is very common to demand an explanation for the numerical score in order for the analysis to be believed before publication or utilization in a court of law, however, most deepfake detection methods lack such an explanation. If results for single interactions between humans and computers are to be accepted it is vitally important to involve XAI early in the process.

### B. Robust PAD for ASV in Smart Homes

Because the Smart Home concept encourages hands-free automation, bio-metric technologies, and particularly automatic speech recognition, are the ideal way to regulate and personalize access. However, while contemporary spoofing systems are extremely descriptive and specific when detecting presentation attacks, current speaker identification systems may be vulnerable. Therefore, most of these PAD systems are ineffective in the current environment where the types of presentation attacks are unknown. When deployed to smart homes, the risk of attack increases significantly. Hence, there is an immediate need to build a strong PAD system to protect ASV systems in order to accelerate the safe use of ASVs in Smart Home applications.

### C. What type of model is best for given scenario?

Because countermeasures operate differently in varying conditions there is an urgent need to shift the research community's focus to environmental or scenario-based spoofing countermeasures. The environment is an essential part of spoofing detection countermeasures; for instance, several countermeasures may fail to provide adequate performance when surrounded by a loud environment or background noise. As a result, a thorough comparison study is required to find the best countermeasure for specific scenarios.

### D. Need for a new Audio Forensics datasets

Voice-controlled devices (VCDs) have transformed the Internet of Things (IoT), enabling the enhanced development, personalization, and automation, of smart homes through voice activation. Unfortunately, these VCDs are currently largely unprotected, and can be exploited via spoofing attacks. Existing datasets, such as the Voice spoofing detection corpus (VSDC), ASVspoof 2017, ASVspoof 2019, and ReMASC, contain large-scale replay (i.e., single and multi-order replays) and basic cloning audio files, but newer, deep learning-based attacks may easily defeat current solutions, introducing a new level of complexity for anti-spoofing algorithms. Moreover, large-scale datasets like ASVspoof 2017 and ASVspoof 2019 lack cloning design parameters, which are an essential component of current VCDs. Therefore, an enhanced cloning dataset is needed in order to stay ahead of the curve of advanced voice cloning spoof attacks. For creation of such dataset, new generative algorithms such as diffusion models could be employed. Also, there is need to develop datasets that contains multiple forgeries at frame level in a single file. The new datasets also need to be unbiased with respect to gender and ethnicity. subsectionThe role of Adversarial Machine Learning on existing Audio spoofing countermeasures High-performance spoofing countermeasure systems for automated speaker verification (ASV) have been presented during the ASVspoof sessions. However, the robustness of these systems in the face of adversarial attacks has yet to be thoroughly examined, while adversarial attacks pose a significant threat to ASV systems and their countermeasures. Currently, there are only a limited number of countermeasures in this domain. It is therefore necessary to develop a more adversarial-aware ASV system.

### E. Fairness-enabled voice spoofing countermeasures and ASV

Fairness is defined as the lack of bias or preference toward a person or a group based on intrinsic or acquired attributes. Prior to their widespread commercialization, it is necessary to study the bias and impartiality of spoofing detection and speaker verification systems across populations. No quantitative evaluation has been reported to study the fairness (w.r.t gender, ethnicity etc) in the spoofing detectors across protected parameters, and even integrated voice countermeasure. To avoid the real-world consequences of a biased and flawed system that favors a single sub-group or class, it is vital to ensure the impartiality of data sets used for training and testing spoofing systems.

### F. Attack-aware Audio Forensics Countermeasures

ASV (automated speaker verification) systems should always discard both non-target (voice produced by unrelated speakers) and spoofed (reprocessed or transformed) signals. However, little consideration has been given to how ASV systems should be changed when they encounter spoofing attacks, or even when they pair up with spoofing countermeasures, much less to how both systems might be jointly optimized.

### G. Need for Cross-dataset Evaluation

The need for cross dataset evaluation is important to consider when evaluating the performance of any aspect of an audio spoof detection system as shown in this paper. Cross dataset evaluation gives important insight into the robustness of the tested features and can provide additional information and direction toward a more general and complete ASV system.

### H. Physics of AI

Current deepfake detectors face challenges, particularly due to incomplete, sparse, and noisy data in training phases. There is a need to explore innovative AI architectures, algorithms, and approaches that "bake in" the physics, mathematics, and prior knowledge relevant to deep fakes and other forgeries. Embedding physics and prior knowledge, using knowledge-infused learning, into AI will help to overcome the challenges of sparse data and will facilitate the development of generative models that are causal and explanatory.

## XI. Conclusion

This research provides a comprehensive review of voice spoofing attacks, countermeasures proposed to counter certain attacks, publicly accessible datasets, and the criteria used to evaluate countermeasure performance in voice spoofing detection. This paper also includes an experimental comparison of the state-of-the-art (SOTA) voice spoofing countermeasures against multi-and cross-datasets and classifiers. Based on the results of the experiments, this study concludes that characteristics that better capture microphone signatures and harmonic distortions give improved detection performance for the identification of PA attacks. Furthermore, we show that the efficacy of the SOTA countermeasure varies significantly with changes in datasets and classifiers. Extensive testing reveals that these countermeasures, along with CNN and CNN-GRU classifiers, perform well. Detailed experiments also demonstrate the significance of cross-corpus evaluation for future voice spoofing technologies.

## XII. Reproducibility

We have provided a GitHub repository[2] that contains all the code and explanation needed to reproduce results of models evaluated to test their generalizability. The README.pdf file contains explanation of how to run the experiments.

[2]https://github.com/smileslab/Comparative-Analysis-Voice-Spoofing

# XIII. Acknowledgement

# References

[1] Federico Alegre, Asmaa Amehraye, and Nicholas Evans. Spoofing countermeasures to protect automatic speaker verification from voice conversion. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 3068–3072. IEEE, 2013.

[2] PR Aravind, Usamath Nechiyil, Nandakumar Paramparambath, et al. Audio spoofing verification using deep convolutional neural networks by transfer learning. *arXiv preprint arXiv:2008.03464*, 2020.

[3] Tuba Arif, Ali Javed, Mohammed Alhameed, Fathe Jeribi, and Ali Tahir. Voice spoofing countermeasure for logical access attacks detection. *IEEE Access*, 9:162857–162868, 2021.

[4] B. T. Balamurali, Kinwah Edward Lin, Simon Lui, Jer-Ming Chen, and Dorien Herremans. Toward robust audio spoofing detection: A detailed comparison of traditional and learned features. *IEEE Access*, 7:84229–84241, 2019.

[5] BT Balamurali, Kinwah Edward Lin, Simon Lui, Jer-Ming Chen, and Dorien Herremans. Toward robust audio spoofing detection: A detailed comparison of traditional and learned features. *IEEE Access*, 7:84229–84241, 2019.

[6] Roland Baumann, Khalid Mahmood Malik, Ali Javed, Andersen Ball, Brandon Kujawa, and Hafiz Malik. Voice spoofing detection corpus for single and multi-order audio replays. *Computer Speech & Language*, 65:101132, 2021.

[7] Stephen R Bond, Anke Hoeffler, and Jonathan RW Temple. Gmm estimation of empirical growth models. *Available at SSRN 290522*, 2001.

[8] Weicheng Cai, Danwei Cai, Wenbo Liu, Gang Li, and Ming Li. Countermeasures for automatic speaker verification replay spoofing attack: On data augmentation, feature representation, classification and fusion. In *INTERSPEECH*, pages 17–21, 2017.

[9] Sandipan Chakraborty and Goutam Saha. Improved text-independent speaker identification using fused mfcc & imfcc feature sets based on gaussian filter. *International Journal of Signal Processing*, 5(1):11–19, 2009.

[10] Nanxin Chen, Yanmin Qian, Heinrich Dinkel, Bo Chen, and Kai Yu. Robust deep feature for spoofing detection—the sjtu system for asvspoof 2015 challenge. In *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.

[11] Xinhui Chen, You Zhang, Ge Zhu, and Zhiyao Duan. UR Channel-Robust Synthetic Speech Detection System for ASVspoof 2021. In *Proc. 2021 Edition of the Automatic Speaker Verification and Spoofing Countermeasures Challenge*, pages 75–82, 2021.

[12] Zhuxin Chen, Zhifeng Xie, Weibin Zhang, and Xiangmin Xu. ResNet and Model Fusion for Automatic Spoofing Detection. In *Proc. Interspeech 2017*, pages 102–106, 2017.

[13] Bhusan Chettri, Tomi Kinnunen, and Emmanouil Benetos. Deep generative variational autoencoding for replay spoof detection in automatic speaker verification, 2020.

[14] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. Voxceleb2: Deep speaker recognition. *arXiv preprint arXiv:1806.05622*, 2018.

[15] Héctor Delgado, Nicholas Evans, Tomi Kinnunen, Kong Aik Lee, Xuechen Liu, Andreas Nautsch, Jose Patino, Md Sahidullah, Massimiliano Todisco, Xin Wang, et al. Asvspoof 2021: Automatic speaker verification spoofing and countermeasures challenge evaluation plan. *arXiv preprint arXiv:2109.00535*, 2021.

[16] Héctor Delgado, Massimiliano Todisco, Md Sahidullah, Nicholas Evans, Tomi Kinnunen, Kong Aik Lee, and Junichi Yamagishi. Asvspoof 2017 version 2.0: meta-data analysis and baseline enhancements. In *Odyssey 2018-The Speaker and Language Recognition Workshop*, 2018.

[17] H. Dinkel, N. Chen, Y. Qian, and K. Yu. End-to-end spoofing detection with raw waveform cldnns. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4860–4864, 2017.

[18] Serife Kucur Ergünay, Elie Khoury, Alexandros Lazaridis, and Sebastien Marcel. On the vulnerability of speaker verification to realistic voice spoofing. In *2015 IEEE 7th International Conference on Biometrics Theory, Applications and Systems (BTAS)*, pages 1–6. IEEE, 2015.

[19] Roberto Font, Juan M Espín, and María José Cano. Experimental analysis of features for replay attack detection-results on the asvspoof 2017 challenge. In *Interspeech*, pages 7–11, 2017.

[20] Julie Gerlings, Arisa Shollo, and Ioanna Constantiou. Reviewing the need for explainable artificial intelligence (xai), 2020.

[21] Yuan Gong, Jian Yang, Jacob Huber, Mitchell MacKnight, and Christian Poellabauer. Remasc: Realistic replay attack corpus for voice controlled systems. *arXiv preprint arXiv:1904.03365*, 2019.

[22] Tharshini Gunendradasan, Saad Irtza, Eliathamby Ambikairajah, and Julien Epps. Transmission line cochlear model based am-fm features for replay attack detection. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6136–6140, 2019.

[23] Tharshini Gunendradasan, Buddhi Wickramasinghe, Phu Le, Eliathamby Ambikairajah, and Julien Epps. Detection of replay-spoofing attacks using frequency modulation features. In *INTERSPEECH*, pages 636–640, 09 2018.

[24] Farman Hassan and Ali Javed. Voice spoofing countermeasure for synthetic speech detection. In *2021 International Conference on Artificial Intelligence (ICAI)*, pages 209–212, 2021.

[25] Lian Huang and Chi-Man Pun. Audio replay spoof attack detection by joint segment-based linear filter bank feature extraction and attention-enhanced densenet-bilstm network. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:1813–1825, 2020.

[26] Ting Huang, Hongxia Wang, Yi Chen, and Peisong He. Gru-svm model for synthetic speech detection. In Hongxia Wang, Xianfeng Zhao, Yunqing Shi, Hyoung Joong Kim, and Alessandro Piva, editors, *Digital Forensics and Watermarking*, pages 115–125, Cham, 2020. Springer International Publishing.

[27] Arindam Jati, Chin-Cheng Hsu, Monisankha Pal, Raghuveer Peri, Wael AbdAlmageed, and Shrikanth Narayanan. Adversarial attack and defense strategies for deep speaker recognition systems. *Computer Speech & Language*, 68:101199, Jul 2021.

[28] Ali Javed, Khalid Mahmood Malik, Aun Irtaza, and Hafiz Malik. Towards protecting cyber-physical and iot systems from single-and multi-order voice spoofing attacks. *Applied Acoustics*, 183:108283, 2021.

[29] Ali Javed, Khalid Mahmood Malik, Hafiz Malik, and Aun Irtaza. Voice spoofing detector: A unified anti-spoofing framework. *Expert Systems with Applications*, 198:116770, 2022.

[30] Zhe Ji, Zhi-Yi Li, Peng Li, Maobo An, Shengxiang Gao, Dan Wu, and Faru Zhao. Ensemble learning for countermeasure of audio replay spoofing attack in asvspoof2017. In *Interspeech*, pages 87–91, 2017.

[31] Akhil Jose, Justin Joseph, Glan Devadhas, and M. M. Shinu. Influence of filter bank structure on the statistical significance of coefficients in cepstral analysis for acoustic signals. In Sabu M. Thampi, S ri Krishnan, Juan Manuel Corchado Rodriguez, Swagatam Das, Michal Wozniak, and Dhiya Al-Jumeily, editors, *Advances in Signal Processing and Intelligent Recognition Systems*, pages 91–104, Cham, 2018. Springer International Publishing.

[32] 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, 2021.

[33] 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 *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6367–6371. IEEE, 2022.

[34] Jee-weon Jung, Hye-jin Shim, Hee-Soo Heo, and Ha-Jin Yu. Replay attack detection with complementary high-resolution information using end-to-end dnn for the asvspoof 2019 challenge. *arXiv preprint arXiv:1904.10134*, 2019.

[35] Jee-weon Jung, Hemlata Tak, Hye-jin Shim, Hee-Soo Heo, Bong-Jin Lee, Soo-Whan Chung, Ha-Jin Yu, Nicholas Evans, and Tomi Kinnunen. Sasv 2022: The first spoofing-aware speaker verification challenge. *arXiv preprint arXiv:2203.14732*, 2022.

[36] Madhu R Kamble, Hardik B Sailor, Hemant A Patil, and Haizhou Li. Advances in anti-spoofing: from the perspective of asvspoof challenges. *APSIPA Transactions on Signal and Information Processing*, 9, 2020.

[37] Tomi Kinnunen, Kong Aik Lee, Héctor Delgado, Nicholas Evans, Massimiliano Todisco, Md Sahidullah, Junichi Yamagishi, and Douglas A Reynolds. t-dcf: a detection cost function for the tandem assessment of spoofing countermeasures and automatic speaker verification. *arXiv preprint arXiv:1804.09618*, 2018.

[38] Tomi Kinnunen, Md Sahidullah, Héctor Delgado, Massimiliano Todisco, Nicholas Evans, Junichi Yamagishi, and Kong Aik Lee. The asvspoof 2017 challenge: Assessing the limits of replay spoofing attack detection. 2017.

[39] Tomi Kinnunen, Md Sahidullah, Héctor Delgado, Massimiliano Todisco, Nicholas Evans, Junichi Yamagishi, and Kong Aik Lee. The asvspoof 2017 challenge: Assessing the limits of replay spoofing attack detection. In *INTERSPEECH*, 08 2017.

[40] Tomi H. Kinnunen, Md. Sahidullah, Mauro Falcone, Luca Costantini, Rosa González Hautamäki, Dennis Alexander Lehmann Thomsen, Achintya Kumar Sarkar, Z. Tan, Héctor Delgado, Massimiliano Todisco, Nicholas W. D. Evans, Ville Hautamäki, and Kong-Aik Lee. Reddots replayed: A new replay spoofing attack corpus for text-dependent speaker verification research. *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5395–5399, 2017.

[41] Pavel Korshunov, André R Gonçalves, Ricardo PV Violato, Flávio O Simões, and Sébastien Marcel. On the use of convolutional neural networks for speech presentation attack detection. In *2018 IEEE 4th international conference on identity, security, and behavior analysis (ISBA)*, pages 1–8. IEEE, 2018.

[42] Pavel Korshunov, Sébastien Marcel, Hannah Muckenhirn, André R Gonçalves, AG Souza Mello, RP Velloso Violato, Flávio O Simoes, M Uliani Neto, Marcus de Assis Angeloni, José Augusto Stuchi, et al. Overview of btas 2016 speaker anti-spoofing competition. In *2016 IEEE 8th international conference on biometrics theory, applications and systems (BTAS)*, pages 1–6. IEEE, 2016.

[43] Felix Kreuk, Yossi Adi, Moustapha Cisse, and Joseph Keshet. Fooling end-to-end speaker verification with adversarial examples. In *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 1962–1966. IEEE, 2018.

[44] Jia Min Karen Kua, Tharmarajah Thiruvaran, Mohaddeseh Nosratighods, Eliathamby Ambikairajah, and Julien Epps. Investigation of spectral centroid magnitude and frequency for speaker recognition. In *Odyssey*, page 7, 2010.

[45] Cheng-I Lai, Alberto Abad, Korin Richmond, Junichi Yamagishi, Najim Dehak, and Simon King. Attentive filtering networks for audio replay attack detection. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6316–6320. IEEE, 2019.

[46] Cheng-I Lai, Nanxin Chen, Jesús Villalba, and Najim Dehak. Assert: Anti-spoofing with squeeze-excitation and residual networks. *arXiv preprint arXiv:1904.01120*, 2019.

[47] Galina Lavrentyeva, Sergey Novoselov, Andzhukaev Tseren, Marina Volkova, Artem Gorlanov, and Alexandr Kozlov. Stc antispoofing systems for the asvspoof2019 challenge. *arXiv preprint arXiv:1904.05576*, 2019.

[48] Kong-Aik Lee, Anthony Larcher, Guangsen Wang, Patrick Kenny, Niko Brümmer, David A. van Leeuwen, Hagai Aronowitz, Marcel Kockmann, Carlos Vaquero, Bin Ma, Haizhou Li, Themos Stafylakis, Md. Jahangir Alam, Albert Swart, and Javier Pérez. The reddots data collection for speaker recognition. In *INTERSPEECH*, 2015.

[49] Xu Li, Na Li, Chao Weng, Xunying Liu, Dan Su, Dong Yu, and Helen Meng. Replay and synthetic speech detection with res2net architecture. In *ICASSP 2021-2021 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 6354–6358. IEEE, 2021.

[50] Leian Liu and Jichen Yang. Study on feature complementarity of statistics, energy, and principal information for spoofing detection. *IEEE Access*, 8:141170–141181, 2020.

[51] Xuechen Liu, Md Sahidullah, and Tomi Kinnunen. Spoofing-aware speaker verification with unsupervised domain adaptation. *arXiv preprint arXiv:2203.10992*, 2022.

[52] Youxuan Ma, Zongze Ren, and Shugong Xu. Rw-resnet: A novel speech anti-spoofing model using raw waveform, 2021.

[53] K. M. Malik, A. Javed, H. Malik, and A. Irtaza. A light-weight replay detection framework for voice controlled iot devices. *IEEE Journal of Selected Topics in Signal Processing*, 14(5):982–996, 2020.

[54] Khalid Mahmood Malik, Ali Javed, Hafiz Malik, and Aun Irtaza. A light-weight replay detection framework for voice controlled iot

devices. *IEEE Journal of Selected Topics in Signal Processing*, 14(5):982–996, 2020.

[55] Khalid Mahmood Malik, Ali Javed, Hafiz Malik, and Aun Irtaza. A light-weight replay detection framework for voice controlled iot devices. *IEEE Journal of Selected Topics in Signal Processing*, 14(5):982–996, 2020.

[56] Aakshi Mittal and Mohit Dua. Automatic speaker verification systems and spoof detection techniques: review and analysis. *International Journal of Speech Technology*, pages 1–30, 2021.

[57] João Monteiro, Jahangir Alam, and Tiago H Falk. An ensemble based approach for generalized detection of spoofing attacks to automatic speaker recognizers. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6599–6603. IEEE, 2020.

[58] Lindasalwa Muda, Mumtaj Begam, and Irraivan Elamvazuthi. Voice recognition algorithms using mel frequency cepstral coefficient (mfcc) and dynamic time warping (dtw) techniques. *arXiv preprint arXiv:1003.4083*, 2010.

[59] Hema A Murthy and Venkata Gadde. The modified group delay function and its application to phoneme recognition. In *2003 IEEE International Conference on Acoustics, Speech, and Signal Processing, 2003. Proceedings.(ICASSP'03).*, volume 1, pages I–68. IEEE, 2003.

[60] Parav Nagarsheth, Elie el Khoury, Kailash Patil, and Matt Garland. Replay attack detection using dnn for channel discrimination. In *INTERSPEECH*, 2017.

[61] Arsha Nagrani, Joon Son Chung, and Andrew Zisserman. Voxceleb: a large-scale speaker identification dataset. *arXiv preprint arXiv:1706.08612*, 2017.

[62] Ravika Naika. An overview of automatic speaker verification system. In *Intelligent Computing and Information and Communication*, pages 603–610. Springer, 2018.

[63] Sergey Novoselov, Alexandr Kozlov, Galina Lavrentyeva, Konstantin Simonchik, and Vadim Shchemelinin. Stc anti-spoofing systems for the asvspoof 2015 challenge. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5475–5479. IEEE, 2016.

[64] Hemant A Patil and Madhu R Kamble. A survey on replay attack detection for automatic speaker verification (asv) system. In *2018 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pages 1047–1053. IEEE, 2018.

[65] Dipjyoti Paul, Monisankha Pal, and Goutam Saha. Novel speech features for improved detection of spoofing attacks. In *2015 Annual IEEE India Conference (INDICON)*, pages 1–6. IEEE, 2015.

[66] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

[67] Padmanabhan Rajan, Tomi Kinnunen, Cemal Hanilci, Jouni Pohjalainen, and Paavo Alku. Using group delay functions from all-pole models for speaker recognition. In *INTERSPEECH*, pages 2489–2493. Citeseer, 2013.

[68] Amir Mohammad Rostami, Mohammad Mehdi Homayounpour, and Ahmad Nickabadi. Efficient attention branch network with combined loss function for automatic speaker verification spoof detection. *arXiv preprint arXiv:2109.02051*, 2021.

[69] Md Sahidullah, Héctor Delgado, Massimiliano Todisco, Tomi Kinnunen, Nicholas Evans, Junichi Yamagishi, and Kong-Aik Lee. Introduction to voice presentation attack detection and recent advances. In *Handbook of biometric anti-spoofing*, pages 321–361. Springer, 2019.

[70] Md Sahidullah, Héctor Delgado, Massimiliano Todisco, Hong Yu, Tomi Kinnunen, Nicholas Evans, and Zheng-Hua Tan. Integrated spoofing countermeasures and automatic speaker verification: An evaluation on asvspoof 2015. 2016.

[71] Md Sahidullah, Tomi Kinnunen, and Cemal Hanilçi. A comparison of features for synthetic speech detection, 2015.

[72] M. S. Saranya, R. Padmanabhan, and Hema A. Murthy. Replay attack detection in speaker verification using non-voiced segments and decision level feature switching. In *2018 International Conference on Signal Processing and Communications (SPCOM)*, pages 332–336, 2018.

[73] Ibon Saratxaga, Inmaculada Hernáez, Michael Pucher, and Inaki Sainz. Perceptual importance of the phase related information in speech. In *INTERSPEECH*, volume 2, 09 2012.

[74] Gajan Suthokumar, Kaavya Sriskandaraja, Vidhyasaharan Sethu, Chamith Wijenayake, and Eliathamby Ambikairajah. Phoneme specific modelling and scoring techniques for anti spoofing system. In *ICASSP*

*2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6106–6110, 2019.

[75] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks, 2014.

[76] Hemlata Tak, Jee-weon Jung, Jose Patino, Madhu Kamble, Massimiliano Todisco, and Nicholas Evans. End-to-end spectro-temporal graph attention networks for speaker verification anti-spoofing and speech deepfake detection, 2021.

[77] Hemlata Tak, Jee-weon Jung, Jose Patino, Massimiliano Todisco, and Nicholas Evans. Graph attention networks for anti-spoofing, 2021.

[78] Hemlata Tak, Madhu Kamble, Jose Patino, Massimiliano Todisco, and Nicholas Evans. Rawboost: A raw data boosting and augmentation method applied to automatic speaker verification anti-spoofing, 2021.

[79] Choon Beng Tan, Mohd Hanafi Ahmad Hijazi, Norazlina Khamis, Zuraini Zainol, Frans Coenen, Abdullah Gani, et al. A survey on presentation attack detection for automatic speaker verification systems: State-of-the-art, taxonomy, issues and future direction. *Multimedia Tools and Applications*, 80(21):32725–32762, 2021.

[80] Prasad A Tapkir, Ankur T Patil, Neil Shah, and Hemant A Patil. Novel spectral root cepstral features for replay spoof detection. In *2018 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pages 1945–1950. IEEE, 2018.

[81] Prasad A Tapkir and Hemant A Patil. Significance of teager energy operator phase for replay spoof detection. In *2018 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pages 1951–1956. IEEE, 2018.

[82] Zhongwei Teng, Quchen Fu, Jules White, Maria E Powell, and Douglas C Schmidt. Sa-sasv: An end-to-end spoof-aggregated spoofing-aware speaker verification system. *arXiv preprint arXiv:2203.06517*, 2022.

[83] Massimiliano Todisco, Héctor Delgado, and Nicholas W. D. Evans. Constant q cepstral coefficients: A spoofing countermeasure for automatic speaker verification. *Comput. Speech Lang.*, 45:516–535, 2017.

[84] Massimiliano Todisco, Xin Wang, Ville Vestman, Md Sahidullah, Héctor Delgado, Andreas Nautsch, Junichi Yamagishi, Nicholas Evans, Tomi Kinnunen, and Kong Aik Lee. Asvspoof 2019: Future horizons in spoofed and fake audio detection. *arXiv preprint arXiv:1904.05441*, 2019.

[85] Xavier Valero and Francesc Alias. Gammatone cepstral coefficients: Biologically inspired features for non-speech audio classification. *IEEE Transactions on Multimedia*, 14(6):1684–1689, 2012.

[86] Ville Vestman, Tomi Kinnunen, Rosa González Hautamäki, and Md Sahidullah. Voice mimicry attacks assisted by automatic speaker verification. *Computer Speech & Language*, 59:36–54, 2020.

[87] SVM Vishwanathan and M Narasimha Murty. Ssvm: a simple svm algorithm. In *Proceedings of the 2002 International Joint Conference on Neural Networks. IJCNN'02 (Cat. No. 02CH37290)*, volume 3, pages 2393–2398. IEEE, 2002.

[88] Qian Wang, Xiu Lin, Man Zhou, Yanjiao Chen, Cong Wang, Qi Li, and Xiangyang Luo. Voicepop: A pop noise based anti-spoofing system for voice authentication on smartphones. In *IEEE INFOCOM 2019 - IEEE Conference on Computer Communications*, pages 2062–2070, 2019.

[89] 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*, 64:101114, 2020.

[90] Marcin Witkowski, Stanislaw Kacprzak, Piotr Zelasko, Konrad Kowalczyk, and Jakub Galka. Audio replay attack detection using high-frequency features. In *Interspeech*, pages 27–31, 2017.

[91] Eddie Wong and Sridha Sridharan. Comparison of linear prediction cepstrum coefficients and mel-frequency cepstrum coefficients for language identification. In *Proceedings of 2001 International Symposium on Intelligent Multimedia, Video and Speech Processing. ISIMP 2001 (IEEE Cat. No. 01EX489)*, pages 95–98. IEEE, 2001.

[92] Haibin Wu, Po chun Hsu, Ji Gao, Shanshan Zhang, Shen Huang, Jian Kang, Zhiyong Wu, Helen Meng, and Hung yi Lee. Spotting adversarial samples for speaker verification by neural vocoders, 2021.

[93] Haibin Wu, Xu Li, Andy T. Liu, Zhiyong Wu, Helen Meng, and Hung yi Lee. Adversarial defense for automatic speaker verification by cascaded self-supervised learning models, 2021.

[94] Haibin Wu, Xu Li, Andy T. Liu, Zhiyong Wu, Helen Meng, and Hung yi Lee. Improving the adversarial robustness for speaker verification by self-supervised learning, 2021.

[95] Haibin Wu, Andy T. Liu, and Hung yi Lee. Defense for black-box attacks on anti-spoofing models by self-supervised learning, 2020.

[96] Haibin Wu, Songxiang Liu, Helen Meng, and Hung yi Lee. Defense against adversarial attacks on spoofing countermeasures of asv, 2020.

[97] Xiang Wu, Ran He, Zhenan Sun, and Tieniu Tan. A light cnn for deep face representation with noisy labels. *IEEE Transactions on Information Forensics and Security*, 13(11):2884–2896, 2018.

[98] Zhenzong Wu, Rohan Kumar Das, Jichen Yang, and Haizhou Li. Light convolutional neural network with feature genuinization for detection of synthetic speech attacks, 2020.

[99] Zhizheng Wu, Nicholas Evans, Tomi Kinnunen, Junichi Yamagishi, Federico Alegre, and Haizhou Li. Spoofing and countermeasures for speaker verification: A survey. *speech communication*, 66:130–153, 2015.

[100] Zhizheng Wu, Sheng Gao, Eng Siong Cling, and Haizhou Li. A study on replay attack and anti-spoofing for text-dependent speaker verification. In *Signal and Information Processing Association Annual Summit and Conference (APSIPA), 2014 Asia-Pacific*, pages 1–5, 2014.

[101] Zhizheng Wu, Ali Khodabakhsh, Cenk Demiroglu, Junichi Yamagishi, Daisuke Saito, Tomoki Toda, and Simon King. Sas: A speaker verification spoofing database containing diverse attacks. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 4440–4444. IEEE, 2015.

[102] Zhizheng Wu, Tomi Kinnunen, Nicholas Evans, Junichi Yamagishi, Cemal Hanilçi, Md Sahidullah, and Aleksandr Sizov. Asvspoof 2015: the first automatic speaker verification spoofing and countermeasures challenge. In *Sixteenth annual conference of the international speech communication association*, 2015.

[103] Zhizheng Wu, Junichi Yamagishi, Tomi Kinnunen, Cemal Hanilçi, Mohammed Sahidullah, Aleksandr Sizov, Nicholas Evans, Massimiliano Todisco, and Héctor Delgado. Asvspoof: the automatic speaker verification spoofing and countermeasures challenge. *IEEE Journal of Selected Topics in Signal Processing*, 11(4):588–604, 2017.

[104] Zhizheng Wu, Junichi Yamagishi, Tomi Kinnunen, Cemal Hanilçi, Mohammed Sahidullah, Aleksandr Sizov, Nicholas Evans, Massimiliano Todisco, and Héctor Delgado. Asvspoof: The automatic speaker verification spoofing and countermeasures challenge. *IEEE Journal of Selected Topics in Signal Processing*, 11(4):588–604, 2017.

[105] Jichen Yang and Rohan Das. Long-term high frequency features for synthetic speech detection. *Digital Signal Processing*, 97:102622, 11 2019.

[106] Jichen Yang and Rohan Kumar Das. Low frequency frame-wise normalization over constant-q transform for playback speech detection. *Digital Signal Processing*, 2019.

[107] Jichen Yang and Rohan Kumar Das. Improving anti-spoofing with octave spectrum and short-term spectral statistics information. *Applied Acoustics*, 157:107017, 2020.

[108] Jichen Yang, Rohan Kumar Das, and Nina Zhou. Extraction of octave spectra information for spoofing attack detection. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 27(12):2373–2384, 2019.

[109] Jichen Yang, Longting Xu, Bo Ren, and Yunyun Ji. Discriminative features based on modified log magnitude spectrum for playback speech detection, 2020.

[110] You Zhang, Fei Jiang, and Zhiyao Duan. One-class learning towards synthetic voice spoofing detection. *IEEE Signal Processing Letters*, 28:937–941, 2021.

[111] You Zhang, Ge Zhu, and Zhiyao Duan. A probabilistic fusion framework for spoofing aware speaker verification. *arXiv preprint arXiv:2202.05253*, 2022.

[112] Fang Zheng, Guoliang Zhang, and Zhanjiang Song. Comparison of different implementations of mfcc. *Journal of Computer science and Technology*, 16(6):582–589, 2001.

[113] Xinhui Zhou, Daniel Garcia-Romero, Ramani Duraiswami, Carol Espy-Wilson, and Shihab Shamma. Linear versus mel frequency cepstral coefficients for speaker recognition. In *2011 IEEE Workshop on Automatic Speech Recognition & Understanding*, pages 559–564. IEEE, 2011.