

# Multilingual Audio Deepfakes Dataset for Robust and Generalizable Detection

Candy Olivia Mawalim\*, Yutong Wang\*, Aulia Adila\*, Shogo Okada, Masashi Unoki

*Graduate School of Advanced Science and Technology*

*Japan Advanced Institute of Science and Technology, Ishikawa, Japan*

{candylin, s2310404, adila, okada-s, unoki}@jaist.ac.jp

**Abstract**—The rise of sophisticated technologies capable of generating realistic synthetic human speech has introduced significant security challenges in voice-based applications. These advances in speech generation have enabled malicious actors to produce convincing audio deepfakes, undermining the reliability of speaker verification systems and digital communication. Despite growing interest in audio deepfake detection, many existing datasets are limited in linguistic diversity and fail to capture the complexity of real-world scenarios, thus constraining model generalization. In this work, we introduce the JMAD (JAIST Multilingual Audio Deepfake) dataset, a large-scale, multilingual, and multi-source corpus designed to support the development and evaluation of robust and generalizable audio deepfake detection systems. Covering 17 languages and comprising over 400,000 utterances, JMAD incorporates a wide range of deepfake generation methods and includes both human (pristine) and machine-generated speech, sourced from publicly available corpora and a curated subset of private data. We provide detailed analyses of utterance duration, generation techniques, and audio quality, along with comprehensive evaluations across multiple model architectures and configurations. Cross-dataset evaluations are also conducted to assess the generalization capabilities of detection models across diverse languages and data domains. This study contributes to a deeper understanding of the limitations and opportunities in current detection systems, ultimately paving the way for more resilient and linguistically inclusive countermeasures.

**Index Terms**—audio deepfake detection, multilingual dataset, generalizability

## I. INTRODUCTION

The advancement of deep learning technologies has significantly improved the quality of generated content across various modalities. Among these, audio synthesis technologies such as text-to-speech (TTS) and voice conversion (VC) have enabled a wide range of beneficial applications, particularly in human communication. However, they have also introduced serious potential threats to social security and political stability when misused for malicious purposes [1]. One such threat is the use of audio deepfakes, which refers to audio that has been digitally altered using artificial intelligence (AI) [1], with the intent to deceive humans. This can lead to harmful consequences such as fraud and the spread of misinformation. Another concern involves biometric identification systems, where manipulated audio is used to bypass automatic speaker verification (ASV) systems [2], [3].

An emerging defense strategy against audio deepfakes is spoofing detection, which aims to distinguish between human and machine-generated (spoofed) speech. Developing a reliable detection model presents several challenges. First, the continuous evolution of speech synthesis techniques, including diverse model architectures and training algorithms, necessitates that detection models generalize effectively across a wide variety of attacks. This highlights the importance of training detection systems on realistic and representative datasets to ensure robustness. Furthermore, as speech generation technologies expand into multilingual settings [4]–[6], it becomes increasingly critical to advance spoofing detection capabilities beyond high-resource languages in order to ensure inclusive protection.

In response to these challenges, we introduce the JMAD (JAIST Multilingual Audio Deepfake) dataset, a new initiative designed to facilitate the development of robust detection models capable of distinguishing human from machine-generated speech across multilingual and multisource conditions. The JMAD dataset comprises a comprehensive multilingual speech corpus compiled from multiple open-source resources and small parts of internally collected data, rigorously curated to ensure both quality and representativeness. Spanning 17 languages and encompassing a broad spectrum of synthesis methods, recording conditions, and speaker demographics, the JMAD dataset is explicitly constructed to improve model generalization beyond what existing resources currently support.

To assess the utility of our dataset, we present a comprehensive evaluation encompassing several key aspects. First, we evaluate the inherent quality of the dataset. Second, we quantify the detection performance using benchmark methods specifically on the JMAD dataset. Finally, we conduct a cross-dataset evaluation to demonstrate its potential for training robust models capable of generalizing to unseen data.

## II. RELATED WORK

The most widely adopted dataset for advancing the development of audio deepfake detection systems is provided by the ASVspoof Challenge series [3], [7], [8], [24], [25], which has served as a benchmark primarily for English. Earlier editions focused on spoofing attacks targeting automatic speaker verification (ASV) systems, while more recent editions have expanded to include standalone countermeasures independent of ASV. The latest edition, ASVspoof 5 [25], is built upon the

\*These authors contributed equally.

TABLE I: Existing datasets for audio deepfake detection, a comparison based on the language coverage, year, audio deepfake generation methods, number of utterances (all, pristine, and generated), as well as the primary focus task addressed in each dataset

Dataset	Year	Language(s)	Generation methods	# Utts.	# Pristine	# Generated	Focus task(s)
ASVspoof2015 [7]	2015	English	TTS, Vocoder, VC	263,151	16,651	246,500	Spoofing detection (speech synthesis and voice conversion)
ASVspoof2017 [8]	2017	English	Replay	18,030	3,565	14,465	Spoofing detection (replay attacks)
Fake or Real [9]	2019	English	TTS	>198,000	>110,000	>87,000	General audio deepfake detection
ASVspoof2019 [3]	2019	English	TTS, Vocoder, VC, Replay	100,448	10,256	90,192	Spoofing detection (speech synthesis, voice conversion, and replay attacks)
WaveFake [10]	2021	English, Japanese	TTS, Vocoder	117,985	0	117,985	General audio deepfake detection
ASVspoof2021 [11]	2021	English	TTS, Vocoder, VC, Replay	593,253	20,637	572,616	General audio deepfake and spoofing detection (speech synthesis, voice conversion, replay attacks)
ADD2022 [12]	2022	Mandarin	TTS, VC, Partially Fake	160,885	36,953	123,932	General audio deepfake detection
CFAD [13]	2022	Mandarin	TTS	347,400	115,800	231,600	General audio deepfake detection (robustness, generalization)
In-the-Wild [14]	2022	English	TTS	31,779	19,963	11,816	Real-world audio deepfake detection
ADD2023 [15]	2023	Mandarin	TTS, VC, Partially Fake	195,541	172,819	113,042	General audio deepfake detection (include manipulation region location and deepfake algorithm recognition)
DEEP-VOICE [16]	2023	English	VC (retrieval-based)	7,484	3,742	3,742	Voice conversion deepfakes detection
DECRO [17]	2023	English, Mandarin	TTS, VC	118,381	33,702	84,679	Cross-language audio deepfake detection
MLAAD [18]	2024	Multilingual (38 lang.)	TTS, Vocoder	154,000	0	154,000	General audio deepfake detection
CD-ADD [19]	2024	English	TTS	145,570	25,111	120,459	Cross-domain audio deepfake detection
ASVspoof5 [20]	2024	English	TTS, Vocoder, VC, AT	1,500,713	289,527	1,211,186	General audio deepfake and spoofing detection (including adversarial attacks)
DFADD [21]	2024	English	TTS	207,955	44,455	163,500	Spoofing detection (diffusion and flow-matching based TTS)
CVoiceFake [22]	2024	Multilingual (5 lang.)	Vocoder	1,254,893	0	1,254,893	Multilingual audio deepfake detection and speech content privacy preservation
SAFE Challenge [23]	2025	Multilingual	Unknown	Unknown	Unknown	Unknown	Unseen and various audio deepfake detection (robustness)
<b>JMAD-Open (ours)</b>	2025	Multilingual (15 lang.)	TTS, Vocoder, VC, AT	312,146	65,546	246,600	Multilingual audio deepfake detection
<b>JMAD-All (ours)</b>	2025	Multilingual (17 lang.)	TTS, Vocoder, VC, AT	412,021	93,968	318,053	Multilingual audio deepfake detection

MLS [26] corpus and incorporates adversarial attacks applied to spoofed utterances generated using various TTS, vocoder, and VC algorithms. It also introduces codec simulation to reflect real-world audio transmission scenarios.

Another well-known initiative is the ADD Challenge [12], [15], which addresses more complex real-world detection scenarios. The first edition focused on detecting low-quality and partially fake audio [12], while the second edition expanded to include manipulation localization and generation method identification [15]. The ADD dataset is derived from Mandarin corpora, AISHELL-1 [27], AISHELL-3 [28], and AISHELL-4 [29], and contains samples generated using a range of TTS and VC models, though the specific models are not publicly disclosed.

Filling the gap, CFAD [13] was introduced as the public Mandarin standard dataset for fake audio detection under noisy and transcoding conditions. It consists of three dataset versions: clean, noisy, and codec. Human speech samples were sourced from both open datasets and self-recorded data (six sources in total). It incorporates 12 type of fake audios, 11 of which are generated using synthesis methods with different vocoders, and the remaining one is partially fake obtained by

clipping and splicing.

Additionally, several English-language corpora have been proposed to advance research in audio deepfake detection. The Fake or Real (FoR) dataset [9], introduced in 2019, includes samples generated using a combination of open-source and commercial tools. The In-the-Wild dataset [14], designed to evaluate model generalization in real-world conditions, comprises found audio recordings of celebrities and politicians, with approximately half of the recordings being deepfakes. Another notable resource is the DFADD dataset [21], which contains deepfake audio generated using five state-of-the-art diffusion and flow-matching TTS models. To address the growing challenges posed by zero-shot TTS systems, the Cross-Domain Audio Deepfake Detection (CD-ADD) dataset [19] was developed to support detection models across varying domains. A separate dataset was also introduced utilizing retrieval-based voice conversion systems such as DEEP-VOICE [16]. Furthermore, to support cross-lingual evaluation of detection systems, the DECRO dataset [17] was created, incorporating speech samples in both English and Mandarin.

Subsequently, several datasets have been developed to support languages beyond English and Chinese. The WaveFake

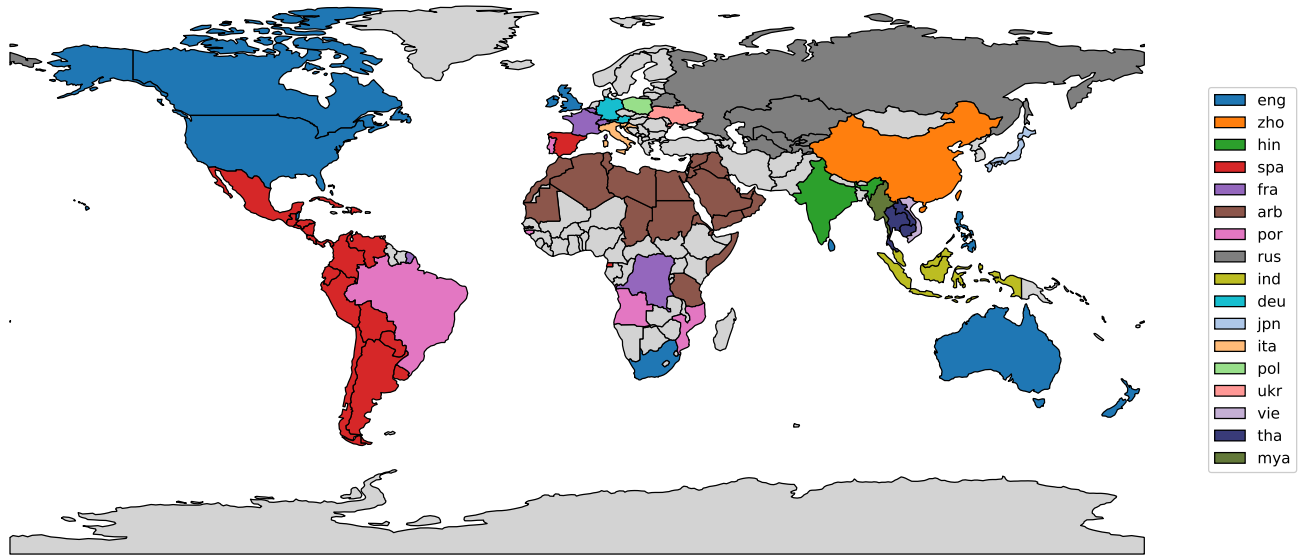


Fig. 1: World map illustrating the primary geographical regions where the 17 languages included in our multilingual corpus are spoken as de facto and/or de jure languages. The legend provides the corresponding ISO 639 language codes for each region.

dataset [10] includes generated audio synthesized using neural vocoder models, covering both English and Japanese. A prominent multilingual dataset is MLAAD [18], which spans 38 languages and is built upon the M-AILABS dataset [30], originally composed of recordings in eight languages sourced from audiobooks and interviews. The generated speech in MLAAD is synthesized using 82 TTS models across 33 architectures, including the Griffin-Lim vocoder. MLAAD has demonstrated strong utility by enabling the training of deepfake detection models that outperform those trained on other datasets such as In-the-Wild and Fake or Real. Its extensive linguistic diversity also facilitates robust cross-lingual evaluation and improves generalization. Another large-scale multilingual dataset is CVoiceFake [22], which contains over 1.25 million bonafide and deepfake utterances across five languages, with speech data sourced from the CommonVoice dataset [31].

Recent efforts have also focused on detecting audio deepfakes in low-resource languages. Notably, studies have targeted languages within the ASEAN region [32], aiming to develop spoofing countermeasures tailored to these linguistic contexts. This line of work includes the construction of dedicated datasets for Thai (ThaiSpoof [33]), Indonesian (InaSpoof [34], [35]), Vietnamese (VSASV [36]), and Myanmar (UC-SYSpooF [37]). These studies underscore several challenges in building effective detection models for underrepresented languages—such as limited access to high-quality human speech data, inconsistencies in dataset quality across languages, and the rapid advancement of realistic spoofing techniques.

Table I summarizes the comparison of existing datasets for audio deepfake detection. Although these datasets have

contributed significantly to the field, most are still limited in linguistic diversity and lack balanced representation across languages. Moreover, few incorporate both pristine and generated speech from a wide variety of sources, which is critical for building generalizable and robust detection models.

### III. JMAD (JAIST MULTILINGUAL AUDIO DEEPPFAKE) DATASET

To facilitate the development of a reliable model for detecting generated speech across multiple languages, we assembled and carefully curated a diverse multilingual speech dataset, drawing from various open and private sources to ensure high quality and broad linguistic coverage.

#### A. Curation Process

To begin the dataset curation process, as illustrated in Fig. 2, we first compiled a comprehensive list of publicly available speech corpora spanning various languages. From these, we carefully selected 17 representative spoken languages, including English, Mandarin, Hindi, Italian, Modern Standard Arabic, Spanish, Polish, German, French, Russian, Portuguese, Japanese, Ukrainian, Vietnamese, Thai, Indonesian, and Myanmar. Figure 1 illustrates the geographical areas where these languages function as official or widely spoken primary languages.

To ensure high-fidelity human speech (pristine) and maintain relevance with current detection system technologies, we prioritized widely used public corpora, primarily developed for spoofing detection, speech recognition, or synthesis tasks. Table II provides a detailed breakdown of the dataset composition, including:

TABLE II: Overview of data sources used in the JMAAD dataset, detailing the adopted languages, data types (pristine and/or generated), deepfake methods (TTS/text-to-speech, vocoder, VC/voice conversion, AT/adversarial attack), recording or synthesis conditions, audio format, and sample rate.

Source	Adopted language(s)	Type(s)	Deepfake methods	Audio condition	Format	Sample rate
ASVspoof [20]	English	Pristine, Generated	TTS, Vocoder, VC, AT	Clean with post processing	FLAC	16 kHz
ADD [12]	Mandarin	Generated	TTS, VC	Noisy	WAV	16 kHz
MLAAD [18]	Hindi, Spanish, French, Arabic, Thai Portuguese, Russian, German, Japanese, Italian, Polish, Ukrainian, Vietnamese	Generated	TTS, Vocoder	Clean	WAV	22 kHz
AISHELL-3 [28]	Mandarin	Pristine	-	Clean	WAV	16 kHz
M-AILABS [30]	Spanish, French, Russian, German, Italian, Polish, Ukrainian	Pristine	-	Unknown	WAV	16 kHz
MediaSpeech [38]	Arabic	Pristine	-	Unknown	WAV	16 kHz
Indic-voice [39]	Hindi	Pristine	-	Clean, Noisy	WAV	8 kHz
Private data	Myanmar, Indonesian, Vietnamese English (nonnative), Japanese, Thai	Pristine, Generated	TTS, Vocoder, VC	Clean, Noisy	Multiple	Multiple

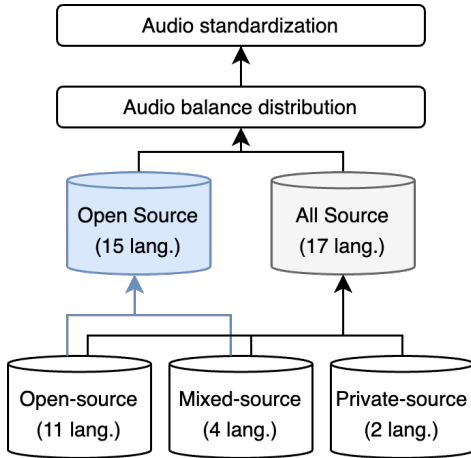


Fig. 2: JMAAD dataset curation process, comprising the audio source selection and acquisition, audio balance distribution, and audio standardization.

- the source repository name(s) and corresponding languages used in our dataset;
- the type of speech included (pristine, generated, or both), along with the associated deepfake or generation algorithms;
- a summary of the speech conditions, such as noise levels, audio formats, and sampling rates.

For English, we used well-established corpora from the ASVspoof Challenge [25] for both pristine and generated speech, which comprise the largest portion of our dataset. We chose the latest challenge version to best align with current detection system development. This version builds upon the MLS English data [26] and includes stronger attacks, featuring advanced text-to-speech (TTS), vocoder, voice conversion (VC), and adversarial attacks (AT) designed to mislead automatic speaker verification (ASV) and countermeasure (CM)

systems. To simulate real-world transmission conditions, codec compression was also applied to both pristine and generated audio to simulate realistic audio transmission conditions.

For Chinese, we adopted corpora from the Audio Deepfake Detection (ADD) Challenge [12] to supplement generated data. The pristine Chinese data was sourced from AISHELL-3 [28], which is the same corpus referenced by the ADD Challenge, ensuring the inclusion of high-quality human speech.

To support multilingual diversity in both pristine and generated data, we incorporated established corpora such as M-AILABS [30] and MLAAD [18]. In addition, language-specific resources like IndicVoice [39] and MediaSpeech [38] were used to supplement pristine data in the corresponding languages, i.e. Hindi and Arabic.

In some cases, open-source corpora lacked sufficient quantity or quality of human speech. Where feasible, we supplemented the dataset with internally recorded speech. For instance, we included English utterances spoken by non-native speakers from Asia to increase accent diversity. These hybrid sources are referred to as ‘mixed-source’ repositories and also apply to languages such as Japanese, Vietnamese, and Thai. The inclusion of these recordings helped ensure broad linguistic coverage and representativeness.

Private-source data were also incorporated from prior studies on spoofing detection in Asian languages [32]–[35], [37], [40], each containing pristine and generated speech synthesized from a variety of TTS, vocoder, and VC algorithms. The pristine recordings were typically collected in diverse environments using multiple devices, enriching the dataset’s acoustic variability and thereby improving generalization for detection models.

We further observed that some repositories, such as MLAAD, included only a single speaker to generate spoofed data. Therefore, we limited the number of utterances to a maximum of 100 per spoofing algorithm in single-speaker settings, while retaining all utterances in multi-speaker settings from

TABLE III: Summary of the number of utterances (#utts.) distribution across 17 languages in the JMAD dataset.

Language	Open Source (#utts.)		All Source (utts.)	
	Pristine	Generated	Pristine	Generated
English (eng)	50,131	175,040	51,931	175,040
Mandarin (zho)	4,410	30,267	4,410	30,267
Hindi (hin)	4,850	1,959	4,850	1,959
Italian (ita)	2,060	4,263	2,060	4,263
Arabic (arb)	2,505	2,957	2,505	2,957
Spanish (spa)	270	4,742	270	4,742
Polish (pol)	100	4,870	100	4,870
German (deu)	476	4,316	476	4,316
French (fra)	298	4,413	298	4,413
Russian (rus)	148	3,711	148	3,711
Portugese (por)	0	3,011	0	3,011
Japanese (jpn)	0	2,978	5,037	2,978
Ukrainian (ukr)	298	2,237	298	2,237
Vietnamese (vie)	0	961	4,968	961
Thai (tha)	0	875	5,263	24,305
Indonesian (ind)	0	0	5,863	27,960
Myanmar (mya)	0	0	5,491	20,063
<b>Total</b>	<b>65,546</b>	<b>246,600</b>	<b>93,968</b>	<b>318,053</b>
<b>Subtotal</b>	<b>312,146</b>		<b>412,021</b>	

other speech corpora. Additionally, for datasets containing pristine speech from multiple speakers, such as M-AILABS, we carefully selected approximately 50 utterances per speaker. These efforts were made to ensure a balanced representation of speakers in each corpus.

After carefully conduct the data acquisition from open-source, mixed-source, and private-source repositories, we organized the dataset into two configurations. The *Open Source* subset comprises only publicly available corpora (both open- and mixed-source) and serves as the primary contribution for promoting transparency and reproducibility. The *All Source* configuration additionally includes a limited amount of internally recorded speech (private-source), expanding language coverage and offering stronger baselines, particularly for languages with limited availability of high-quality public human speech data.

To reduce bias, we aimed for a balanced distribution of speech samples across different languages, setting the number of machine-generated utterances to be roughly three times that of human speech. This ratio reflects practical scenarios often encountered in voice deepfake detection. Additionally, given that many existing detection systems are primarily developed and tested using English data, we intentionally included a larger volume of English samples. An overview of the dataset

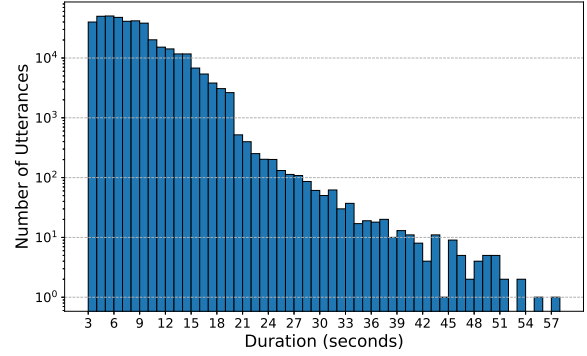


Fig. 3: Histogram showing the distribution of utterance durations in the JMAD dataset (All Source) on a logarithmic scale

composition across 17 languages is presented in Table III.

During the final audio standardization phase, we limited each audio clip to a maximum of 60 seconds to ensure manageability and consistency during training and evaluation. All audio files were resampled to 16 kHz and converted to mono-channel format. We also verified audio integrity by checking for file corruption, empty clips, and ensuring compliance with the .wav format. Conducting such a meticulous data curation process has allowed us to compile a well-structured, multilingual speech corpus that is consistent, diverse, and suitable for advancing research in audio deepfake detection.

### B. Metadata

Our standardized speech corpus is organized into a unified dataset structure consisting of audio samples and their corresponding metadata. Each audio file is annotated with the following metadata fields:

- General information: Identifiers for the utterance, speaker, speaker gender, and a label indicating whether the sample is genuine human speech (labeled as ‘pristine’) or machine-generated (labeled as ‘generated’).
- Codec type: Specifies the audio codec applied to the sample before it is converted to WAV format, if any (e.g., M4A).
- Deepfake algorithm: Describes the method or model used to generate the synthetic or manipulated audio sample.

### C. Statistics

Our compiled dataset comprises a total of 412,021 speech samples across 17 languages, including English and Chinese, which represent some of the world’s most widely spoken languages.

To analyze the temporal characteristics of the dataset, we examined the distribution of utterance durations, as depicted in Fig. 3. Most samples fall within the range of approximately 3 to 10 seconds. Furthermore, we present a breakdown of utterance duration distributions across different languages for both pristine and generated data, as shown in Fig. 7. The



distributions are shown on a logarithmic scale to accommodate the wide variation in utterance counts across languages. Most languages exhibit a reasonable spread across duration categories, which supports generalization for models trained to detect audio artifacts across diverse speaking styles and utterance lengths, thereby enhancing the robustness of evaluation systems under real-world conditions.

Our compiled dataset encompasses a multitude of deepfake methods for speech generation, categorized into four distinct types: TTS-based attacks, vocoder-based attacks, VC-based attacks, and adversarial attacks (AT). TTS-based attacks are generated using text-to-speech systems that synthesize speech directly from textual input, often leveraging models such as GlowTTS, VITS, Tacotron2, or other pretrained neural architectures. Vocoder-based attacks generate speech from intermediate acoustic features such as mel spectrograms, using vocoders like Griffin-Lim. VC-based attacks modify a source speaker’s voice to resemble that of a target speaker, typically without changing the linguistic content, using models such as StarGANv2-VC or ASR-based VC pipelines. Meanwhile, AT introduces subtle, imperceptible perturbations to utterances, specifically optimized to degrade the performance of spoofing detection systems [20].

We incorporated generated speech samples from four different data sources: ASVspoof [25], ADD [12], MLAAD [18], and a portion of our private-source data. Among all data sources included in the JMAD dataset, only ASVspoof explicitly incorporates adversarial attacks as part of its deepfake generation strategy. In this case, adversarial perturbations are applied as a post-processing step to spoofed utterances produced by TTS, vocoder, or VC systems. Portions of both pristine and generated samples have also been subjected to audio encoding and compression using codecs, simulating real-world conditions where speech may be transmitted over networks or stored in compressed formats.

For the MLAAD set, we selected 10 speech generation algorithms, along with their variants and modifications, including TTS models and vocoders tailored to the adopted languages in the JMAD dataset. The ADD dataset also contains generated samples produced using commonly used TTS and VC algorithms, although the exact models are not disclosed. Our private-source generated speech was created using high-performance speech synthesis systems, encompassing TTS, vocoders, VC, and several proprietary TTS models.

A comprehensive list of the generation algorithms included in the JMAD dataset, along with their corresponding data sources, is presented in Table VIII in the Appendix. The term ‘system’ was chosen to reflect the inclusive scope of the deepfake generation methods, encompassing models (e.g., VITS, Tacotron2), techniques or components (e.g., Griffin-Lim vocoder, Malafide attack), as well as frameworks or pipelines (e.g., Whisper-based TTS, unit-selection-based TTS). Additionally, for clarity, we grouped certain distinct algorithms that are built upon the same base model under a single category.

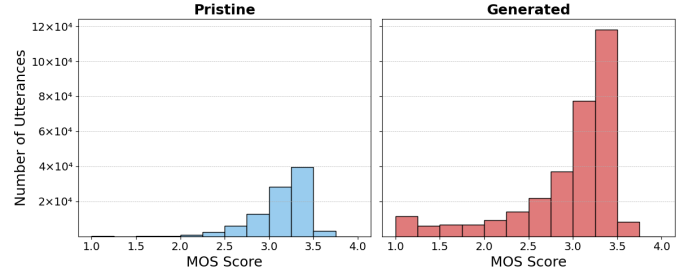


Fig. 4: Histogram showing the distribution of MOS score in the JMAD dataset (All Source)

#### D. Audio Quality

We analyze the quality of the audio data using the Mean Opinion Score (MOS) obtained from DNSMOS [41], a robust and non-intrusive perceptual objective speech quality metric. DNSMOS is well-suited for measuring audio quality as it serves as a proxy for subjective human evaluation, which is considered the “gold standard” in assessing speech quality optimized for human perception. The MOS scale ranges from very poor (MOS = 1) to excellent (MOS = 5). We evaluated both pristine and generated speech across multiple languages and data sources to provide a comprehensive understanding of the dataset’s perceptual quality. A detailed breakdown of the average MOS scores per language is summarized in Table IX in the Appendix, while a comparative view of MOS score distributions across data sources is illustrated in Figure 8 in the Appendix.

Figure 4 displays the distribution of MOS scores for pristine (left) and generated (right) speech in the *All Source* configuration of our dataset. Pristine speech samples generally exhibit moderately high MOS scores, with an average of 3.14, indicating overall good perceptual quality. Only few utterances are scored below 2.5, suggesting that most human recordings are relatively clean and clear. While many generated utterances achieve MOS scores comparable to those of pristine samples, particularly in the 3.0 to 3.5 range, a notable portion falls below 2.5. This drop indicates inconsistencies in quality introduced by certain generation algorithms. Moreover, the broader spread and higher variance in MOS scores for generated speech highlight the varying quality among synthetic samples.

Although this variability in audio quality may help promote generalizability in spoof detection systems, extremely low-quality synthetic speech could introduce undesirable noise during model training. Nevertheless, retaining such low-quality samples in the evaluation set is beneficial, as it ensures the inclusion of various audio conditions and contributes to a more robust final evaluation.

### IV. GENERAL EVALUATION

#### A. Evaluation Metric

Given the straightforward nature of generated audio detection as a binary classification task—distinguishing between

pristine (positive) and generated (negative) audio—we employed balanced accuracy as our evaluation metric. Balanced accuracy, calculated as the average of the accuracy on the pristine class and the accuracy on the generated class, provides a more robust measure of performance, especially in potentially imbalanced datasets. The calculation for balanced accuracy is as follows:

$$\text{Accuracy}_{\text{Pristine}} = \frac{TP}{TP + FP} \quad (1)$$

$$\text{Accuracy}_{\text{Generated}} = \frac{TN}{TN + FN} \quad (2)$$

$$\text{Accuracy}_{\text{Balanced}} = \frac{\text{Accuracy}_{\text{Pristine}} + \text{Accuracy}_{\text{Generated}}}{2} \quad (3)$$

In evaluation, true positive ( $TP$ ) is a correctly identified pristine sample, false positive ( $FP$ ) is a pristine sample incorrectly labeled as generated, true negative ( $TN$ ) is a correctly identified generated sample, and false negative ( $FN$ ) is a generated sample incorrectly labeled as pristine.

In benchmark challenges, Equal Error Rate (EER) and minimum Detection Cost Function (minDCF) often serve as important metrics for audio deepfake detection [42]. EER is the point where the False Acceptance Rate (FAR) and the False Rejection Rate (FRR) of a system are equal. A lower EER generally indicates a more balanced and accurate system, as it signifies a threshold where the trade-off between incorrectly accepting a generated sample as pristine and incorrectly rejecting a pristine sample as generated is minimized.

On the other hand, the minimum Detection Cost Function (minDCF) is a more application-aware metric. It considers the costs associated with both false positives and false negatives, as well as the prior probability of the target class. By minimizing this cost function over different operating points (thresholds), minDCF provides a measure of the best possible performance a system can achieve under specific operational conditions and cost assumptions.

During model development, we also utilized the Area Under the Receiver Operating Characteristic Curve (AUC) to determine an optimal decision threshold for detection. The AUC provides a measure of the model’s ability to distinguish between the pristine and generated classes across various threshold settings, allowing us to select a threshold that balances precision and recall.

## B. Data Partition

We validate our dataset design and collection through experiments in building a robust and generalized spoofing detection.

We partitioned the dataset into three subsets—training, development, and evaluation—using a 6:2:2 ratio. The splitting was primarily based on speaker ID and the source dataset. However, we also considered balancing the distribution of attack ID and gender where this information is available. This process ensures that each subset maintains a similar distribution of these attributes. Furthermore, the three subsets

TABLE IV: Distribution of the JMAD dataset into training (Train), development (Dev), and evaluation (Eval) sets, as used for model validation during our experiments.

Partition	Subset	#utts.		Total
		Pristine	Generated	
All	Train	44,442	137,723	182,165
	Dev	24,977	90,852	115,829
	Eval	24,547	89,478	114,025
Open	Train	28,981	100,452	129,433
	Dev	18,269	73,050	91,319
	Eval	18,296	73,098	91,394

are mutually exclusive, with no overlapping audio samples between them.

Table IV shows the distribution of the training, development, and evaluation sets for all partitioned data used in our experiments. During pre-analysis, we identified a very small number (less than 5) of problematic utterances that caused the total count in the “All” partition to differ from Table III. These problematic utterances have been temporarily removed for the current experiments. In the near future, we plan to investigate the reasons for these issues and implement a fix.

## C. Countermeasures

To evaluate the proposed dataset, we applied several state-of-the-art models. This section outlines the experimental setups, including preprocessing strategies and detailed parameter configurations for each model. We primarily experimented with two model families: Residual Network (ResNet)-based [43] and AASIST-based [44] architectures.

To address the issue of varying audio sample lengths within the dataset, we implemented a preprocessing step to standardize the input duration. Shorter audio clips were padded by repeating their content until the target length was achieved, while longer clips were truncated. Our initial target duration was 4 seconds. However, recognizing that forged audio samples might require more extensive temporal information for accurate classification and to mitigate potential misclassifications due to limited feature representation, we subsequently increased the target durations to 10 seconds. This extension aimed to preserve richer contextual information and ultimately enhance classification accuracy.

1) *CQT-ResNet34*: For our experiments, we employed a ResNet34 model [43] and utilized Constant-Q Transform (CQT) features [45] as input. We conducted experiments using audio segment durations of both 4 seconds and 10 seconds. The CQT features offered a detailed time-frequency representation of the audio, which proved beneficial for capturing the subtle characteristics crucial for distinguishing between genuine and spoofed audio.

2) *RawSpeech-AASIST*: As an end-to-end model, AASIST possesses the capability to directly learn features from raw waveforms [44]. Consequently, for the RawSpeech-AASIST model, we utilized raw waveform inputs directly. Our observations indicated that longer audio segments tended to yield improved model performance, albeit with a corresponding

TABLE V: Cross-partition evaluation on the JMAD dataset, comparing models trained on the OpenSource subset versus the All Source subset. All training configurations used a 4-second audio padding.

			Evaluation data	JMAD-Open				JMAD-All			
Front-end	Back-end	Training data	Label	Accuracy (%)	Balanced Acc. (%)	AUC (%)	EER (%)	Accuracy (%)	Balanced Acc. (%)	AUC (%)	EER (%)
CQT	ResNet34	JMAD-Open	Pristine	0.00	39.99	23.21	69.69	0.00	39.24	24.64	68.10
	Generated		79.98	78.47							
RawSpeech	AASIST		Pristine	59.32	79.15	93.51	12.39	54.27	76.23	91.49	15.39
			Generated	98.97				98.19			
XLS-R (300M)	AASIST		Pristine	<b>81.52</b>	<b>90.66</b>	<b>97.86</b>	<b>4.67</b>	<b>74.88</b>	<b>85.95</b>	<b>97.02</b>	<b>7.06</b>
			Generated	<b>99.79</b>				<b>97.02</b>			
CQT	ResNet34	JMAD-All	Pristine	0.00	39.99	24.62	68.49	0.00	39.24	20.60	71.24
	Generated		79.98	78.47							
RawSpeech	AASIST		Pristine	60.17	79.78	93.94	11.96	66.02	82.56	94.83	10.44
			Generated	99.40				99.09			
XLS-R (300M)	AASIST		Pristine	79.76	89.77	97.59	5.02	73.51	85.83	96.41	7.98
			Generated	99.78				98.15			

increase in computational cost during training and inference. To strike a balance between performance and efficiency, we conducted experiments with input audio durations of 4 and 10 seconds.

3) *SSL-AASIST*: We further investigated a variation of AASIST that integrates self-supervised learning (SSL) features. In this approach, we leveraged pre-trained SSL models to extract rich representations from the raw audio, which were subsequently used as input to the AASIST architecture. Our experiments explored three different input lengths: 4 and 10 seconds. For SSL feature extraction, we selected two high-performing pre-trained models known for their effectiveness across various speech tasks: XLS-R with 300 million parameters **XLS-R 300M** and **WavLM-Large**. These models were chosen for their ability to capture subtle and informative acoustic patterns.

#### D. Results and Analysis

Combining multiple data sources with varying quality can be challenging and may inadvertently degrade the performance of a trained model. To address this, we conducted a thorough cross-evaluation using the two partitions of our dataset: one comprising only open-source data (Open) and the other incorporating private data sources (All). Specifically, we trained representative methods on the Open partition (JMAD-Open) and evaluated them on both the Open and All partitions. Conversely, we also trained on the All partition (JMAD-All) and evaluated on both. The results of this cross-evaluation between the OpenSource and All Source partitions of JMAD are presented in Table V. All method pairs in this evaluation were trained using a default audio padding of 4 seconds.

In our initial cross-partition evaluation on the JMAD dataset, we observed a significant performance disparity across different model architectures. The classical approach employing CQT features with a ResNet34 backbone struggled to accurately detect pristine signals, often classifying all audio as generated, resulting in low balanced accuracy and AUC, and a high EER. In contrast, the end-to-end AASIST architecture demonstrated superior performance. Utilizing raw waveforms (RawSpeech), it achieved a balanced accuracy of approximately 79.15%, an AUC of 93.51%, and an EER of 12.39% when both training and evaluation were performed on JMAD-

Open. Notably, the integration of SSL features, particularly with XLS-R (300M), further enhanced the results, yielding a balanced accuracy of 90.66%, an AUC of 97.86%, and an EER of 4.67%. The performance trends remained consistent when comparing different front-end and back-end model pairings.

Furthermore, our analysis of the Open and All partitions revealed that training on JMAD-Open generally led to better performance compared to training on JMAD-All. This could be attributed to the inclusion of controlled-environment recordings and a low-resource language within the private data of the ‘All’ partition, which increases the difficulty of the detection task. Consequently, the JMAD-All dataset appears to be a more rigorous benchmark for evaluating the robustness of audio deepfake detection models.

In our subsequent analysis, we investigated the impact of different padding lengths on models built using the AASIST architecture. Our experiments revealed that the chosen padding duration is often crucial for accurately detecting artifacts in generated speech, with its importance depending on the characteristics of the evaluation set. For example, a shorter padding might suffice if the evaluation data primarily consists of short utterances (under 5 seconds). Given the diverse utterance lengths within our JMAD dataset, we compared models trained with 4-second and 10-second padding. Figure 5 illustrates the performance of these models, evaluated in terms of balanced accuracy, AUC, EER, and minDCF, on both JMAD-Open and JMAD-All, with training performed on JMAD-Open. Notably, when using raw waveforms as input, the difference in performance between the two padding lengths was minimal. However, we observed a slight improvement when employing XLS-R features with the longer 10-second padding. To further evaluate the generalization capabilities of our models on entirely unseen data, we assessed their performance on the SAFE challenge 2025<sup>1</sup> [23]. Our submissions, utilizing the methods compared in this study, demonstrated a significant performance improvement (over 5% in balanced accuracy) with increased padding size when tested on this completely unknown dataset, which included audio samples with varying lengths up to 60 seconds. A more detailed analysis of our cross-dataset evaluation is presented in Section V.

<sup>1</sup><https://stresearch.github.io/SAFE/>



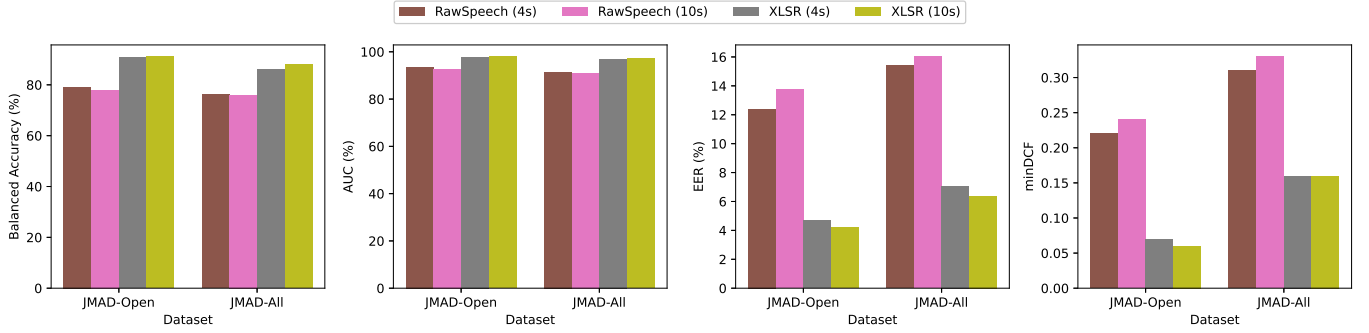


Fig. 5: Performance comparison of AASIST-based models with 4-second and 10-second padding.

TABLE VI: Comparison on various pair of front-end and back-end methods using 10 seconds padding and evaluation on JMAD-Open evaluation data.

Front-end	Back-end	Label	Accuracy (%)	Balanced Acc. (%)	AUC (%)	EER (%)	minDCF
CQT	ResNet	Pristine	25.03	54.19	58.58	44.09	1.00
		Generated	83.34				
XLS-R (300M)	ResNet	Pristine	43.76	68.88	84.96	22.65	0.58
		Generated	94.01				
RawSpeech	AASIST	Pristine	57.17	78.01	92.73	13.75	0.24
		Generated	98.85				
WavLM	AASIST	Pristine	62.12	80.09	95.22	10.99	0.26
		Generated	98.06				
XLS-R (300M)	AASIST	Pristine	82.39	91.14	98.01	4.21	0.06
		Generated	99.88				

To provide a fair comparison of performance across different front-end and back-end method pairings, we trained models on the JMAD-Open dataset using a 10-second padding and evaluated them on the same dataset. The results of this analysis are presented in Table 7, showcasing five combinations: (1) CQT-ResNet, (2) XLSR-ResNet, (3) RawSpeech-AASIST, (4) WavLM-AASIST, and (5) XLSR-AASIST. All evaluation metrics detailed in Subsection IV-A were utilized. Initially, when comparing the traditional CQT features against features extracted from the pre-trained SSL model XLS-R (300M, abbreviated as XLSR), both with a ResNet back-end, we observed a substantial performance difference. Employing the pre-trained SSL features yielded more robust results, albeit at a higher computational cost. Subsequently, the AASIST-based models generally outperformed their ResNet counterparts. For instance, with XLSR as the front-end, AASIST achieved a balanced accuracy of approximately 91.14%, significantly higher than the 78.01% obtained by the ResNet model. Furthermore, our findings highlight the importance of selecting the appropriate SSL model. While WavLM often surpasses XLS-R in various benchmark speech processing tasks [46], in this specific audio deepfake detection scenario, XLSR demonstrated superior performance across all metrics, exhibiting over 10% higher balanced accuracy, a 5% lower EER, and a 0.2 lower minDCF compared to WavLM.

## V. CROSS-DATASET EVALUATION

To assess the generalization capabilities of our models beyond the JMAD dataset, we conducted a comprehensive cross-dataset evaluation. This involved training our models on external datasets, specifically ASVspoof 2019 [3], ASVspoof 2024 [42], and JMAD-Open, and subsequently evaluating their performance on these same datasets as well as the completely unseen SAFE challenge 2025 [23].

Table VII presents the results of our cross-dataset evaluation. It is important to note that due to the lack of ground truth labels for ASVspoof 2024, we were unable to directly evaluate on this benchmark. Instead, we trained models using its training data and assessed their performance on other datasets. Furthermore, the results reported here were achieved with general model architectures and standard hyperparameter tuning on validation sets, without dataset-specific optimizations. The observed performance, particularly in terms of balanced accuracy, reveals considerable variability, highlighting the ongoing challenge of generalization in this domain. While models often perform well when training and evaluation sets are aligned, their effectiveness diminishes on unseen datasets. Consistent with earlier findings, AASIST-based models generally outperformed ResNet.

Notably, our JMAD-Open dataset appears to be a more effective training source for generalization compared to the English-centric ASVspoof 2019 and ASVspoof 2024 datasets.

TABLE VII: Cross-dataset evaluation results for various front-end/back-end method pairings. Training data: ASVspoof2019, ASVspoof2024, JMAD-Open. Evaluation data: ASVspoof2019, JMAD-Open, SAFE challenge 2025 (Task 1).

Front-end	Back-end	Evaluation Data →		ASVspoof2019		JMAD-Open		SAFE (Task1)	
		Training Data ↓	Label	Accuracy (%)	Balanced Acc. (%)	Accuracy (%)	Balanced Acc. (%)	Accuracy (%)	Balanced Acc. (%)
CQT	ResNet	ASVspoof2019	Pristine	35.84	67.40	22.86	53.74	93.15	48.58
			Generated	98.97		84.62		4.00	
RawSpeech	AASIST	ASVspoof2019	Pristine	87.83	93.73	28.85	59.25	4.40	42.70
			Generated	99.64		89.65		81.00	
CQT	ResNet	ASVspoof2024	Pristine	23.17	60.90	26.13	56.76	30.45	36.73
			Generated	98.62		87.40		43.00	
RawSpeech	AASIST	ASVspoof2024	Pristine	27.68	63.20	42.84	68.71	16.90	47.85
			Generated	98.72		94.58		78.79	
CQT	ResNet	JMAD-Open	Pristine	22.86	53.74	25.03	54.19	62.70	58.74
			Generated	84.62		83.34		54.79	
RawSpeech	AASIST		Pristine	27.48	62.94	57.17	78.01	66.10	67.27
			Generated	98.40		98.85		68.43	
WavLM	AASIST		Pristine	31.42	64.89	62.12	80.09	42.95	49.51
			Generated	98.37		98.06		56.07	
XLS-R (300M)	AASIST		Pristine	47.81	73.39	82.39	91.14	72.45	60.73
			Generated	98.97		99.88		49.00	

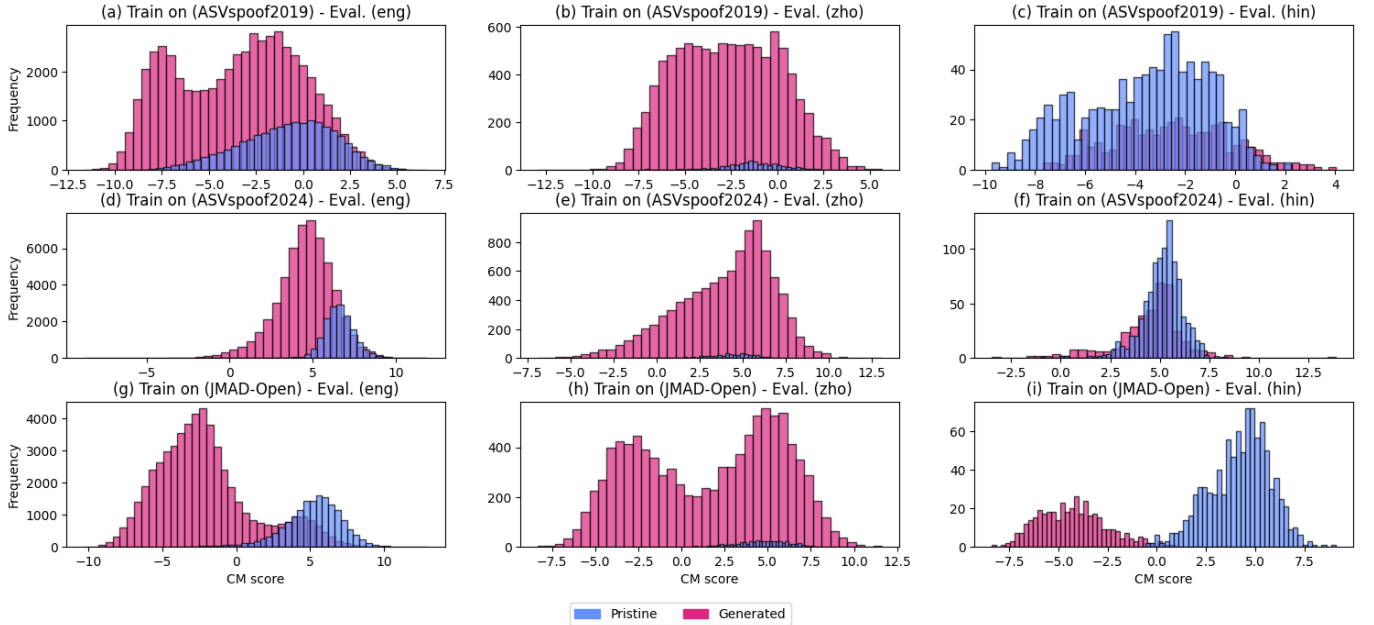


Fig. 6: Language-specific performance using RawSpeech-AASIST method. Three most represented languages within JMAD-Open were included, i.e., English (eng), Mandarin Chinese (zho), and Hindi (hin).

For example, the best performing model trained on JMAD-Open achieved a balanced accuracy of 73.39% on ASVspoof 2019 and 67.27% on the SAFE challenge. Although ASVspoof 2024 showed slightly better performance as training data than ASVspoof 2019, models trained on either ASVspoof dataset struggled significantly on the SAFE challenge, achieving balanced accuracies below 50%.

Further, to investigate the language-specific performance of our best performing architecture, RawSpeech-AASIST, we analyzed the detection results on the three most represented languages within JMAD-Open, i.e., English (eng), Mandarin Chinese (zho), and Hindi (hin). These models were trained on

a combination of ASVspoof 2019, ASVspoof 2024, and the entirety of JMAD-Open. The results are shown in Figure 6.

The language-specific evaluation uncovered notable patterns in the model's generalization capabilities. As anticipated, performance was generally strongest on the model trained on JMAD-Open. The detection accuracy on Mandarin Chinese (zho) across all models, as depicted in Figure 6, was particularly poor. While the countermeasure scores for pristine Mandarin audio evaluated using a model trained on ASVspoof 2024 and JMAD-Open leaned positive, the distribution of generated Mandarin audio significantly overlapped with that of pristine audio. This subpar detection in Mandarin Chinese

could be attributed to the lower quality of generated speech in the Audio Deepfake Database (ADD) benchmark, where the Mean Opinion Score (MOS) is typically below 3.0 (as illustrated in Figure 8b). This language-specific analysis highlights areas for future improvement in multilingual audio deepfake detection.

Moving forward, our work will encompass more extensive cross-dataset evaluations across a broader spectrum of benchmarks. Furthermore, we anticipate that the insights gained from this analysis will aid in the strategic selection of training and evaluation datasets to foster better generalization in audio deepfake detection models.

## VI. DISCUSSION

This study presented the JMAD dataset, a multilingual resource for malicious audio detection, and evaluated its utility through benchmark comparisons and cross-dataset experiments. Our findings highlight several key aspects.

- 1) First, cross-evaluation within JMAD revealed that the partition containing private data, while potentially more challenging, serves as a rigorous testbed.
- 2) Second, the choice of padding length in AASIST-based models significantly impacts performance, particularly on unseen data, as demonstrated by the SAFE challenge results.
- 3) Third, our comparison of front-end and back-end architectures indicated that while SSL-based features enhance performance, AASIST models generally outperform ResNet. Notably, the effectiveness of specific SSL models, such as the superior performance of XLS-R over WavLM in our experiments, underscores the importance of task-specific feature selection.
- 4) Finally, the language-specific analysis on JMAD-Open revealed variations in detection accuracy across languages, with Mandarin Chinese presenting a particular challenge potentially due to the lower quality of available generated speech in the training data.

These findings collectively emphasize the complexities of building robust and generalizable audio deepfake detection systems capable of handling diverse data distributions and linguistic variations.

Despite these contributions, our study has certain limitations. The JMAD dataset, while multilingual, has an uneven distribution of utterances across languages, which might influence the training and evaluation of truly language-agnostic models. Furthermore, our cross-dataset evaluation, while informative, was constrained by the availability of labeled data in external benchmarks, particularly the absence of labels for direct evaluation on ASVspoof 2024. Finally, the hyperparameter tuning for the evaluated models was kept relatively general to assess baseline generalization, and further optimization tailored to specific datasets could potentially yield improved performance. Future work will aim to address these limitations by expanding the linguistic diversity and balance within JMAD, exploring self-supervised learning techniques for improved

cross-dataset generalization, and conducting more extensive hyperparameter optimization for each evaluation scenario.

## VII. CONCLUSION

In conclusion, this work introduced the JMAD dataset as a valuable multilingual resource for advancing research in malicious audio detection. Our evaluations demonstrated its utility for benchmarking and training robust models, highlighting the impact of factors such as data partitioning, input processing (padding), and model architecture. The cross-dataset analysis underscored the challenges of generalization across diverse datasets and languages, while the language-specific findings pointed to the nuanced performance of models across different linguistic contexts. Ultimately, our study contributes to a deeper understanding of the current state of audio deepfake detection and provides a foundation for future research aimed at developing more universally effective and linguistically aware countermeasures.

## DATA AVAILABILITY STATEMENT

The **JAIST Multilingual Audio Deepfake (JMAD)** dataset that curated for this study includes data from both publicly available and private sources. A list of the publicly available datasets used, along with relevant citations can be found in Section III. Due to the inclusion of proprietary data from collaborative projects, the full dataset cannot be made publicly available. However, aggregated statistics and analyses of the dataset are provided within the paper to support our findings. Researchers interested in replicating our work are encouraged to utilize the described publicly available resources.

## ACKNOWLEDGMENT

This work was partially supported by JSPS KAKENHI (25K21245, 23K18491, and 25H01139). Additionally, we extend our sincere thanks to our collaborators in the ASEAN IVO project titled ‘‘Spoof Detection for Automatic Speaker Verification’’ ([www.nict.go.jp/en/asean\\_ivo](http://www.nict.go.jp/en/asean_ivo)) for their prior collaborative work which contributed to this study.

## REFERENCES

- [1] J. Yi, C. Wang, J. Tao, X. Zhang, C. Y. Zhang, and Y. Zhao, ‘‘Audio deepfake detection: A survey,’’ 2023.
- [2] Z. Wu, N. W. D. Evans, T. Kinnunen, J. Yamagishi, F. Alegre, and H. Li, ‘‘Spoofing and countermeasures for speaker verification: A survey,’’ *Speech Commun.*, vol. 66, pp. 130–153, 2015.
- [3] X. Wang, J. Yamagishi, M. Todisco, H. Delgado, A. Nautsch, N. W. D. Evans, M. Sahidullah, V. Vestman, T. H. Kinnunen, K. A. Lee, L. Juvela, P. Alku, Y.-H. Peng, H.-T. Hwang, Y. Tsao, H.-M. Wang, S. L. Maguer, M. Becker, and Z. Ling, ‘‘ASVspoof 2019: A large-scale public database of synthesized, converted and replayed speech,’’ *Comput. Speech Lang.*, vol. 64, p. 101114, 2020.
- [4] A. Babu, C. Wang, A. Tjandra, K. Lakhotia, Q. Xu, N. Goyal, K. Singh, P. von Platen, Y. Saraf, J. Pino, A. Baevski, A. Conneau, and M. Auli, ‘‘XLS-R: self-supervised cross-lingual speech representation learning at scale,’’ in *23rd Annual Conference of the International Speech Communication Association, Interspeech 2022, September 18-22, 2022*. Incheon, Korea: ISCA, 2022, pp. 2278–2282.
- [5] B. Li, Y. Zhang, T. N. Sainath, Y. Wu, and W. Chan, ‘‘Bytes are all you need: End-to-end multilingual speech recognition and synthesis with bytes,’’ in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2019, May 12-17, 2019*. Brighton, United Kingdom: IEEE, 2019, pp. 5621–5625.

- [6] V. Pratap, A. Tjandra, B. Shi, P. Tomasello, A. Babu, S. Kundu, A. Elkahky, Z. Ni, A. Vyas, M. Fazel-Zarandi, A. Baevski, Y. Adi, X. Zhang, W. Hsu, A. Conneau, and M. Auli, "Scaling speech technology to 1,000+ languages," *J. Mach. Learn. Res.*, vol. 25, pp. 97:1–97:52, 2024. [Online]. Available: <https://jmlr.org/papers/v25/23-1318.html>
- [7] Z. Wu, T. Kinnunen, N. W. D. Evans, J. Yamagishi, C. Hanilci, M. Sahidullah, and A. Sizov, "ASVspoof 2015: the first automatic speaker verification spoofing and countermeasures challenge," in *Proc. of INTERSPEECH 2015, September 6-10, 2015*. Dresden, Germany: ISCA, 2015, pp. 2037–2041.
- [8] H. Delgado, M. Todisco, M. Sahidullah, N. W. D. Evans, T. Kinnunen, K. Lee, and J. Yamagishi, "ASVspoof 2017 Version 2.0: meta-data analysis and baseline enhancements," in *Odyssey 2018: The Speaker and Language Recognition Workshop, 26-29 June 2018*. Les Sables d'Olonne, France: ISCA, 2018, pp. 296–303.
- [9] R. Reimao and V. Tzerpos, "For: A dataset for synthetic speech detection," in *2019 International Conference on Speech Technology and Human-Computer Dialogue, SpeD 2019, October 10-12, 2019*. Timisoara, Romania: IEEE, 2019, pp. 1–10.
- [10] J. Frank and L. Schönherr, "WaveFake: A Data Set to Facilitate Audio Deepfake Detection," in *Proc. of NeurIPS Track on Datasets and Benchmarks 2021, December 2021, virtual*, 2021. [Online]. Available: <https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/c74d97b01eae257e44aa9d5bade97baf-Abstract-round2.html>
- [11] J. Yamagishi, X. Wang, M. Todisco, M. Sahidullah, J. Patino, A. Nautsch, X. Liu, K. A. Lee, T. Kinnunen, N. W. D. Evans, and H. Delgado, "ASVspoof 2021: accelerating progress in spoofed and deepfake speech detection," 2021. [Online]. Available: <https://arxiv.org/abs/2109.00537>
- [12] J. Yi, R. Fu, J. Tao, S. Nie, H. Ma, C. Wang, T. Wang, Z. Tian, Y. Bai, C. Fan, S. Liang, S. Wang, S. Zhang, X. Yan, L. Xu, Z. Wen, and H. Li, "ADD 2022: the first audio deep synthesis detection challenge," in *Proc. of ICASSP 2022, 23-27 May 2022*. Virtual and Singapore: IEEE, 2022, pp. 9216–9220.
- [13] H. Ma, J. Yi, C. Wang, X. Yan, J. Tao, T. Wang, S. Wang, and R. Fu, "Cfadd: A chinese dataset for fake audio detection," *Speech Communication*, vol. 164, p. 103122, 2024. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0167639324000931>
- [14] N. Müller, P. Czempin, F. Diekmann, A. Froghyar, and K. Böttinger, "Does audio deepfake detection generalize?" in *Interspeech 2022*, 2022, pp. 2783–2787.
- [15] J. Yi, J. Tao, R. Fu, X. Yan, C. Wang, T. Wang, C. Y. Zhang, X. Zhang, Y. Zhao, Y. Ren, L. Xu, J. Zhou, H. Gu, Z. Wen, S. Liang, Z. Lian, S. Nie, and H. Li, "ADD 2023: the second audio deepfake detection challenge," in *Proceedings of the Workshop on Deepfake Audio Detection and Analysis co-located with 32th International Joint Conference on Artificial Intelligence (IJCAI 2023), August 19, 2023*, ser. CEUR Workshop Proceedings, vol. 3597. Macao, China: CEUR-WS.org, 2023, pp. 125–130. [Online]. Available: <https://ceur-ws.org/Vol-3597/paper21.pdf>
- [16] J. J. Bird and A. Lotfi, "Real-time Detection of AI-Generated Speech for DeepFake Voice Conversion," *CoRR*, vol. abs/2308.12734, 2023.
- [17] Z. Ba, Q. Wen, P. Cheng, Y. Wang, F. Lin, L. Lu, and Z. Liu, "Transferring audio deepfake detection capability across languages," in *Proceedings of the ACM Web Conference 2023*, ser. WWW '23. New York, NY, USA: Association for Computing Machinery, 2023, p. 2033–2044.
- [18] N. M. Müller, P. Kawa, W. H. Choong, E. Casanova, E. Gölge, T. Müller, P. Syga, P. Sperl, and K. Böttinger, "Mlaad: The multi-language audio anti-spoofing dataset," 2025. [Online]. Available: <https://arxiv.org/abs/2401.09512>
- [19] Y. Li, M. Zhang, M. Ren, X. Qiao, M. Ma, D. Wei, and H. Yang, "Cross-domain audio deepfake detection: Dataset and analysis," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Y. Al-Onaizan, M. Bansal, and Y.-N. Chen, Eds. Miami, Florida, USA: Association for Computational Linguistics, Nov. 2024, pp. 4977–4983. [Online]. Available: <https://aclanthology.org/2024.emnlp-main.286/>
- [20] X. Wang, H. Delgado, H. Tak, J. weon Jung, H. jin Shim, M. Todisco, I. Kukanov, X. Liu, M. Sahidullah, T. Kinnunen, N. Evans, K. A. Lee, J. Yamagishi, M. Jeong, G. Zhu, Y. Zang, Y. Zhang, S. Maiti, F. Lux, N. Müller, W. Zhang, C. Sun, S. Hou, S. Lyu, S. L. Maguer, C. Gong, H. Guo, L. Chen, and V. Singh, "ASVspoof 5: Design, Collection and Validation of Resources for Spoofing, Deepfake, and Adversarial Attack Detection Using Crowdsourced Speech," 2025. [Online]. Available: <https://arxiv.org/abs/2502.08857>
- [21] J. Du, I.-M. Lin, I.-H. Chiu, X. Chen, H. Wu, W. Ren, Y. Tsao, H.-Y. Lee, and J.-S. R. Jang, "Dfadd: The diffusion and flow-matching based audio deepfake dataset," in *2024 IEEE Spoken Language Technology Workshop (SLT)*, 2024, pp. 921–928.
- [22] X. Li, K. Li, Y. Zheng, C. Yan, X. Ji, and W. Xu, "Safeear: Content privacy-preserving audio deepfake detection," in *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '24. New York, NY, USA: Association for Computing Machinery, 2024, p. 3585–3599. [Online]. Available: <https://doi.org/10.1145/3658644.3670285>
- [23] S. C. Organizers, "SAFE: Synthetic Audio Forensics Evaluation Challenge," accessed: May 11, 2025. [Online]. Available: <https://stresearch.github.io/SAFE/>
- [24] X. Liu, X. Wang, M. Sahidullah, J. Patino, H. Delgado, T. Kinnunen, M. Todisco, J. Yamagishi, N. W. D. Evans, A. Nautsch, and K. A. Lee, "ASVspoof 2021: Towards Spoofed and Deepfake Speech Detection in the Wild," *IEEE ACM Trans. Audio Speech Lang. Process.*, vol. 31, pp. 2507–2522, 2023.
- [25] X. Wang, H. Delgado, H. Tak, J. Jung, H. Shim, M. Todisco, I. Kukanov, X. Liu, M. Sahidullah, T. Kinnunen, N. W. D. Evans, K. A. Lee, and J. Yamagishi, "Asvspoof 5: Crowdsourced speech data, deepfakes, and adversarial attacks at scale," 2024.
- [26] V. Pratap, Q. Xu, A. Sriram, G. Synnaeve, and R. Collobert, "MLS: A Large-Scale Multilingual Dataset for Speech Research," in *Proc. of INTERSPEECH 2020, October 25-29*. Virtual Event, Shanghai, China: ISCA, 2020, pp. 2757–2761.
- [27] H. Bu, J. Du, X. Na, B. Wu, and H. Zheng, "Aishell-1: An open-source mandarin speech corpus and a speech recognition baseline," in *20th Conference of the Oriental Chapter of the International Coordinating Committee on Speech Databases and Speech I/O Systems and Assessment, O-COCOSDA 2017*. South Korea: IEEE, 11 2017, pp. 1–5.
- [28] Y. Shi, H. Bu, X. Xu, S. Zhang, and M. Li, "Aishell-3: A multi-speaker mandarin tts corpus," in *Interspeech 2021*. Czechia: ISCA, 2021, pp. 2756–2760.
- [29] Y. Fu, L. Cheng, S. Lv, Y. Jv, Y. Kong, Z. Chen, Y. Hu, L. Xie, J. Wu, H. Bu, X. Xu, J. Du, and J. Chen, "Aishell-4: An open source dataset for speech enhancement, separation, recognition and speaker diarization in conference scenario," in *Interspeech 2021*. Czechia: ISCA, 2021, pp. 3665–3669.
- [30] Imdata Celeste, "The M-AILABS Speech Dataset," <https://github.com/imdataceleste/m-ailabs-dataset>, 2020, accessed: April 2025.
- [31] R. Ardila, M. Branson, K. Davis, M. Kohler, J. Meyer, M. Henretty, R. Morais, L. Saunders, F. Tyers, and G. Weber, "Common voice: A massively-multilingual speech corpus," in *Proceedings of the Twelfth Language Resources and Evaluation Conference*, N. Calzolari, F. Béchet, P. Blache, K. Choukri, C. Cieri, T. Declerck, S. Goggi, H. Isahara, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odijk, and S. Piperidis, Eds. Marseille, France: European Language Resources Association, May 2020, pp. 4218–4222. [Online]. Available: <https://aclanthology.org/2020.lrec-1.520/>
- [32] C. O. Mawalim, K. Galajit, D. P. Lestari, W. P. Pa, and M. Unoki, "Challenges in Speech Spoofing Countermeasures for Southeast Asian Languages," *ASJ Spring Meeting 2025*, Saitama, Japan, 2025.
- [33] K. Galajit, T. Kosolsriwiwat, M. Unoki, C. O. Mawalim, P. Aimmamee, W. Kongprawechnon, W. P. Pa, A. Chaiwongyen, T. Racharak, S. Boonkla, H. Yassin, and J. Karnjana, "ThaiSpoof: A Database for Spoof Detection in Thai Language," in *Proc. of (ISAI-NLP)*. Thai: IEEE, 2023, pp. 1–6.
- [34] S. A. Arief, C. O. Mawalim, and D. P. Lestari, "Indonesian Speech Anti-Spoofing System: Data Creation and Convolutional Neural Network Models," in *Proc. of ICAICTA 2024*. Singapore: IEEE, 2024, pp. 1–6.
- [35] C. O. Mawalim, S. A. Arief, and D. P. Lestari, "InaSAS: Benchmarking Indonesian Speech Antispoofing Systems," *APSIPA Transactions on Signal and Information Processing*, 2025, accepted.
- [36] V. Hoang, V. Pham, H. Xuan, P. Nhi, P. Dat, and T. Nguyen, "VSASV: a Vietnamese Dataset for Spoofing-Aware Speaker Verification," in *Proc. Interspeech 2024*. Kos Island, Greece: ISCA, 2024.
- [37] H. M. S. Naing, W. P. Pa, A. M. Hlaing, M. A. A. Aung, K. Galajit, and C. O. Mawalim, "UCSYSPOOF: A Myanmar Language Dataset for Voice Spoofing Detection," in *Proc. of OCOOSDA 2024, October 17-19, 2024*. Hsinchu City, Taiwan: IEEE, 2024, pp. 1–5.

- [38] R. Kolobov, O. Okhapkina, O. Omelchishina, A. Platonov, R. Bedyakin, V. Moshkin, D. Menshikov, and N. Mikhaylovskiy, “MediaSpeech: Multilanguage ASR Benchmark and Dataset,” 2021. [Online]. Available: <https://arxiv.org/abs/2103.16193>
- [39] T. Javed, J. Nawale, E. George, S. Joshi, K. Bhogale, D. Mehendale, I. Sethi, A. Ananthanarayanan, H. Faquih, P. Palit, S. Ravishankar, S. Sukumaran, T. Panchagnula, S. Murali, K. Gandhi, A. R. M. M. C. Vajayanthi, K. Karunganni, P. Kumar, and M. Khapra, “IndicVoices: Towards building an Inclusive Multilingual Speech Dataset for Indian Languages,” in *Findings of the ACL 2024*. Bangkok, Thailand: ACL, Aug. 2024, pp. 10 740–10 782.
- [40] A. Adila, C. O. Mawalim, and M. Unoki, “Detecting Spoof Voices in Asian Non-Native Speech: An Indonesian and Thai Case Study,” in *Proc. of APSIPA ASC 2024, December 3-6*. Macau: IEEE, 2024, pp. 1–6.
- [41] C. K. A. Reddy, V. Gopal, and R. Cutler, “DNSMOS: A Non-Intrusive Perceptual Objective Speech Quality Metric to Evaluate Noise Suppressors,” in *Proc. of (ICASSP) 2021*, 2021, pp. 6493–6497.
- [42] H. Delgado, N. Evans, J.-w. Jung, T. Kinnunen, I. Kukanov, K. A. Lee, X. Liu, H.-j. Shim, M. Sahidullah, H. Tak, M. Todisco, X. Wang, and J. Yamagishi, “ASVspoof 5 Evaluation Plan,” ASVspoof consortium, Tech. Rep., 2024. [Online]. Available: <http://www.asvspoof.org/>
- [43] M. Alzantot, Z. Wang, and M. B. Srivastava, “Deep Residual Neural Networks for Audio Spoofing Detection,” in *Proc. of Interspeech 2019, Graz, Austria, September 15-19*. ISCA, 2019, pp. 1078–1082.
- [44] J. Jung, H. Heo, H. Tak, H. Shim, J. S. Chung, B. Lee, H. Yu, and N. W. D. Evans, “AASIST: audio anti-spoofing using integrated spectro-temporal graph attention networks,” in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2022, 23-27 May 2022*. Virtual and Singapore: IEEE, 2022, pp. 6367–6371.
- [45] J. Brown, “Calculation of a Constant Q Spectral Transform,” *Journal of the Acoustical Society of America*, vol. 89, pp. 425–, 01 1991.
- [46] S. Chen, C. Wang, Z. Chen, Y. Wu, S. Liu, Z. Chen, J. Li, N. Kanda, T. Yoshioka, X. Xiao, J. Wu, L. Zhou, S. Ren, Y. Qian, Y. Qian, J. Wu, M. Zeng, X. Yu, and F. Wei, “WavLM: Large-Scale Self-Supervised Pre-Training for Full Stack Speech Processing,” *IEEE J. Sel. Top. Signal Process.*, vol. 16, no. 6, pp. 1505–1518, 2022.

## VIII. APPENDIX

TABLE VIII: Deepfake methods and systems used to generate samples in the JMAD dataset, along with their corresponding data sources. The ADD dataset is not included due to the undisclosed details of its generation methods.

Deepfake method	Systems	Source
TTS	Bark	MLAAD, Private-source
	XTTS	MLAAD, ASVspoof
	MMS-TTS	MLAAD, Private-source
	VITS	MLAAD, ASVspoof
	Tacotron2	MLAAD
	OpenVoice V2	MLAAD
	Glow-TTS	MLAAD, ASVspoof
	Whisper-based TTS	MLAAD
	vixTTS	MLAAD
	Grad-TTS	ASVspoof
	FastPitch	ASVspoof
	ToucanTTS	ASVspoof
	YourTTS	ASVspoof
	ZMM-TTS	ASVspoof
	Unit selection-based	ASVspoof
Vocoder	Proprietary TTS	Private-source
	Griffin-Lim	MLAAD
	WORLD	Private-source
VC	Hifi-GAN	Private-source
	StarGANv2-VC	ASVspoof
	VAE-GAN	ASVspoof
	In-house ASR-based VC	ASVspoof
	DiffVC	ASVspoof
	FreeVC	Private-source
AT	Spectral filtering	Private-source
	Malafide	ASVspoof
	Malacopula	ASVspoof

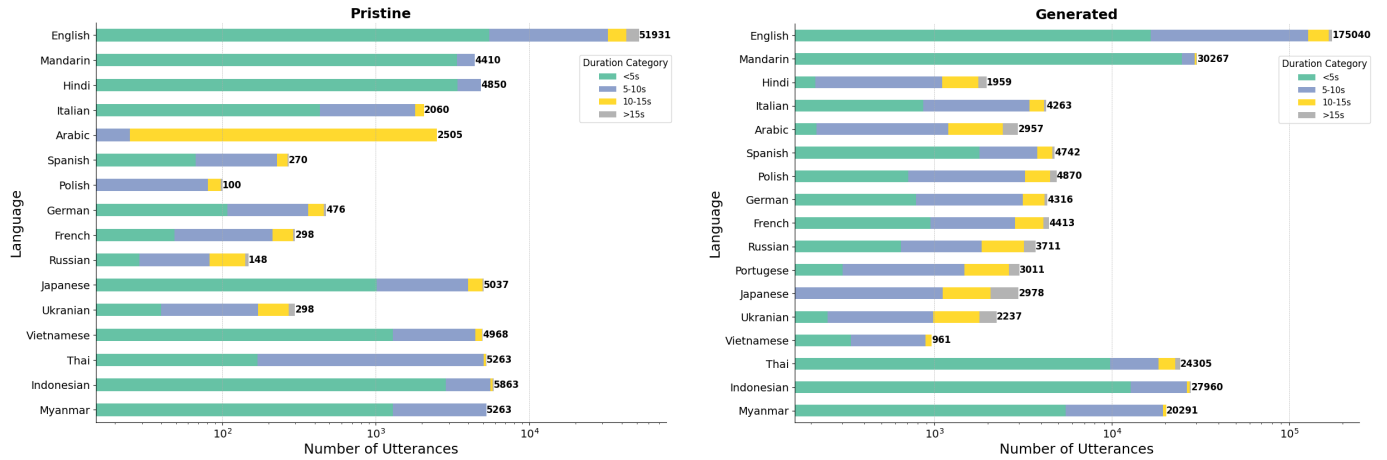


Fig. 7: A detailed breakdown of utterance duration distribution across different languages for pristine (left) and generated (right) samples in the JMAD dataset, illustrated on a logarithmic scale.

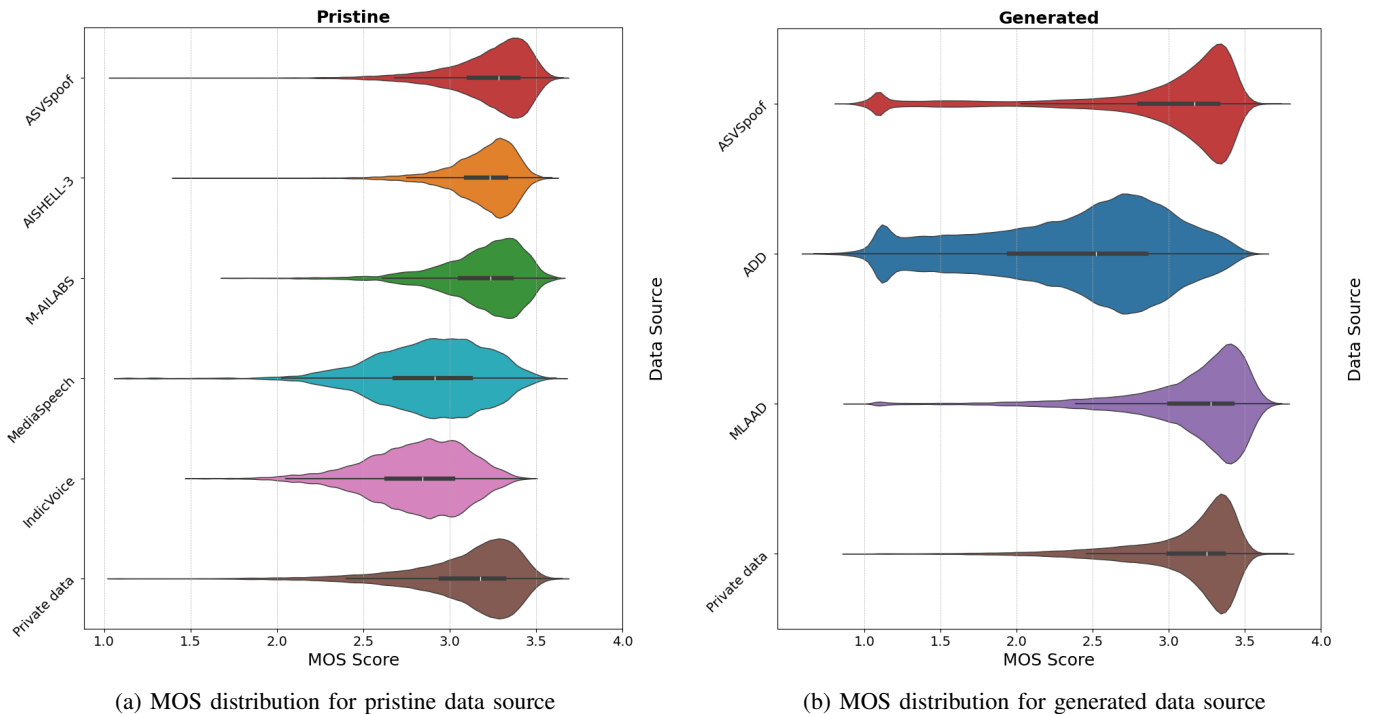


Fig. 8: Violin plots illustrating the distribution of MOS scores for pristine (left) and generated (right) speech across different data sources in the JMAD dataset. Pristine sources generally show consistently high perceptual quality. Among the generated sources, only ADD exhibits a notably broad range of audio quality, indicating varied synthesis characteristics, while others like ASVspoof, MLAAD, and Private data tend to produce more consistent high-MOS outputs.



TABLE IX: Audio quality of the JMAD dataset measured using mean opinion score (MOS) for pristine and generated speech, along with its standard deviation.

Language	Pristine	Generated
English (eng)	$2.91 \pm 0.65$	$3.21 \pm 0.26$
Mandarin (zho)	$2.38 \pm 0.63$	$3.19 \pm 0.20$
Hindi (hin)	$3.14 \pm 0.4$	$2.81 \pm 0.29$
Italian (ita)	$3.04 \pm 0.47$	$3.17 \pm 0.24$
Arabic (arb)	$3.29 \pm 0.31$	$2.88 \pm 0.33$
Spanish (spa)	$3.15 \pm 0.42$	$3.14 \pm 0.28$
Polish (pol)	$3.15 \pm 0.42$	$3.24 \pm 0.24$
German (deu)	$3.09 \pm 0.42$	$3.12 \pm 0.31$
French (fra)	$3.18 \pm 0.38$	$3.26 \pm 0.21$
Russian (rus)	$3.16 \pm 0.41$	$3.29 \pm 0.19$
Portugese (por)	$3.22 \pm 0.4$	
Japanese (jpn)	$3.22 \pm 0.45$	$3.09 \pm 0.35$
Ukranian (ukr)	$2.92 \pm 0.52$	$3.22 \pm 0.22$
Vietnamese (vie)	$3.3 \pm 0.21$	$3.08 \pm 0.33$
Thai (tha)	$3.08 \pm 0.44$	$3.17 \pm 0.25$
Indonesian (ind)	$3.12 \pm 0.38$	$2.95 \pm 0.38$
Myanmar (mya)	$3.14 \pm 0.3$	$3.16 \pm 0.26$
<b>All</b>	<b><math>3.14 \pm 0.3</math></b>	<b><math>2.93 \pm 0.61</math></b>