

Multilingual Dataset Integration Strategies for Robust Audio Deepfake Detection: A SAFE Challenge System

Hashim Ali, Surya Subramani, Lekha Bollinani, Nithin Sai Adupa, Sali El-Loh, Hafiz Malik
Electrical and Computer Engineering
University of Michigan
 Dearborn, USA

Abstract—The SAFE Challenge evaluates synthetic speech detection across three tasks: unmodified audio, processed audio with compression artifacts, and laundered audio designed to evade detection. We systematically explore self-supervised learning (SSL) front-ends, training data compositions, and audio length configurations for robust deepfake detection. Our AASIST-based approach incorporates WavLM large frontend with RawBoost augmentation, trained on a multilingual dataset of 256,600 samples spanning 9 languages and over 70 TTS systems from Codec-Fake, MLAAD v5, SpoofCeleb, Famous Figures, and MAILABS. Through extensive experimentation with different SSL front-ends, three training data versions, and two audio lengths, we achieved second place in both Task 1 (unmodified audio detection) and Task 3 (laundered audio detection), demonstrating strong generalization and robustness.

Index Terms—Synthetic speech detection, audio antispoofing, deepfake detection, AASIST, TTS detection

I. INTRODUCTION

The rapid advancement of text-to-speech (TTS) synthesis technologies has created an urgent need for robust audio deepfake detection systems. Although these technologies offer beneficial applications, their misuse for creating convincing fake audio poses significant security and societal risks. The SAFE (Synthetic Audio Forensics Evaluation) Challenge¹ represents a critical step forward in addressing these challenges by providing a comprehensive evaluation framework that emphasizes real-world applicability and generalizability across diverse audio sources. Unlike traditional evaluation protocols that often focus on controlled laboratory conditions, the SAFE Challenge is designed to test detection systems in realistic scenarios where spoofing attacks may originate from unknown synthesis methods, undergo various processing operations, or be deliberately laundered to evade detection. This emphasis on generalizability across diverse sources reflects the practical challenges faced by deepfake detection systems deployed in real-world environments, where the characteristics of spoofed audio can vary dramatically from training data.

A significant limitation of current deepfake detection research is the reliance on single, often clean datasets for training. Most existing systems are trained on well-curated datasets such as ASVspoof [1]–[3], which, while valuable for controlled evaluation, may not adequately prepare models for

the diversity of spoofing techniques and audio characteristics encountered in practice. This training paradigm can lead to models that perform well on in-domain test sets but fail to generalize to new attack types, different languages, or varying audio quality conditions.

To address this generalization challenge, we conducted four iterative experiments that systematically explored dataset integration strategies for training robust audio deepfake detection systems. Our approach progressively incorporated multiple datasets that span different languages, spoofing techniques, and audio quality conditions, while also investigating the impact of factors such as audio segment length and SSL (self-supervised learning) front-end model selection. These experiments were designed to understand how diverse multilingual datasets can improve the robustness of detection systems when faced with unknown spoofing attacks. We evaluated our approach on both the SAFE Challenge tasks and the In-The-Wild (ITW) dataset [4], a community-standard benchmark to assess generalized performance of audio deepfake detection systems.

Our work makes several practical contributions to the audio deepfake detection community. First, we provide an empirical evaluation of how the combination of datasets affects detection performance across the three SAFE Challenge tasks and on the ITW benchmark [4]. Second, we analyze performance patterns that reveal insights into task-specific challenges, particularly the difficulty of detecting laundered audio. Third, we demonstrate the value of incorporating diverse multilingual training data for improving cross-domain generalization capabilities. Finally, we performed a source-level analysis for both generated and pristine (authentic) sources from the SAFE Challenge, revealing critical vulnerabilities and failure patterns across different synthesis methods and audio processing scenarios. These findings offer actionable insights for researchers and practitioners working to develop more robust audio deepfake detection systems.

The remainder of this paper is organized as follows. Section II reviews existing audio deepfake detection datasets and self-supervised learning approaches relevant to spoofing detection. Section III describes the SAFE Challenge framework and evaluation protocol. Section IV presents our model architecture and SSL front-end selection methodology. Section V details our dataset integration strategies and systematic experimental

¹<https://stresearch.github.io/SAFE/>

design. Section VI covers training configurations and implementation details. Section VII provides comprehensive results and analysis from both SAFE Challenge and ITW evaluations. Finally, Section VIII concludes with implications for the deepfake detection community.

II. RELATED WORK

A. Audio Deepfake Detection Datasets

The development of robust audio deepfake detection systems has been driven by the availability of diverse evaluation datasets, each addressing specific aspects of the spoofing detection challenge.

ASVspoof Series: The ASVspoof Challenges [2], [5]–[7] have provided foundational datasets for the community. ASVspoof 2015 [8] introduced the first spoofing database with 10 TTS systems. ASVspoof 2019 [2] LA expanded this with 19 TTS systems (6 known, 13 unknown), including WaveNet, Tacotron2, and traditional approaches, establishing the standard for controlled evaluation. ASVspoof 2021 [6] introduced two tracks: Logical Access (LA) with 13 systems and Deepfake (DF) with more than 100 attack methods, significantly expanding the attack diversity. The recent ASVspoof 5 [7] represents the most complex design with multilingual data and a wide variety of attacks on a large scale.

Multilingual and Cross-Domain Datasets: MLAAD (Multi-Language Audio Anti-Spoofing Dataset) [9] addresses linguistic diversity with 91 TTS systems across 38 languages, providing 420.7 hours of synthetic speech from 42 different architectures. This dataset specifically targets the language bias present in predominantly English-focused datasets. MLAAD dataset only provides the synthetic audio samples. The M-AILABS Speech Dataset² complements this by providing authentic multilingual speech samples across multiple languages, sourced from audiobooks and public figures speeches.

Real-World and Noisy Conditions: SpoofCeleb [10] leverages VoxCeleb1 [11] as source data, training 23 contemporary TTS systems in real-world noisy conditions to bridge the gap between clean laboratory data and practical applications. In-The-Wild (ITW) [4] provides audio deepfakes collected from social media platforms, offering genuine “in-the-wild” evaluation conditions. Hashim et al. [12] applied various real-world processing on ASVspoof19 audio to generate laundered version of the dataset.

Technology-Specific Datasets: CodecFake [13] introduces the first dataset focused specifically on neural codec-based synthesis, featuring 15 codec models from 6 frameworks, including SpeechTokenizer [14], Encodec [15], and novel approaches like FunCodec [16]. The DFADD dataset [17] focuses on diffusion- and flow-matching-based TTS systems, including GradTTS [18], NaturalSpeech2 [19], Style-TTS2 [20], Matcha-TTS [21], PFlow-TTS [22], etc.

Famous Figures Dataset [23]: Motivated by recent incidents such as the Biden robocall [24], and fabricated recordings of London Mayor Sadiq Khan making inflammatory

remarks [25], we curated a specialized dataset for protecting famous figures from voice cloning attacks. This dataset was developed by collecting high-quality bonafide speech samples of famous figures from YouTube, and then generating the corresponding synthetic speech using various TTS approaches. This dataset provides a deep coverage of political personalities using cutting-edge 2024-2025 TTS systems, such as StyleTTS2 [20], XTTSv2 [26], F5TTS [27], E2TTS [28], etc.

B. Self-Supervised Learning in Audio Spoofing Detection

Self-supervised learning (SSL) has emerged as a powerful approach in audio deepfake detection, addressing challenges such as limited labeled data, generalization to unseen attacks, and robustness across domains. SSL leverages large amounts of unlabeled audio to learn discriminative representations, which can then be fine-tuned or used directly for downstream deepfake detection tasks. Tak et al. [29] investigated the use of Wav2Vec2 for audio spoof detection, demonstrating that SSL-based models achieved SOTA performance even when trained exclusively on bona fide speech samples. They also introduced RawBoost, a data augmentation framework to enhance robustness against real-world distortions.

The cross-dataset generalization capabilities of SSL models have been particularly important. Pascu et al. [30] showed that using frozen SSL representations with simple classifiers significantly improved performance, reducing the Equal Error Rate (EER) from 30.9% to 8.8% across eight deepfake datasets, while emphasizing the importance of model calibration for reliable confidence scores. Recent works [31], [32] have demonstrated that different SSL models capture complementary spoofing artifacts, with WavLM consistently outperforming other models. Stourbe et al. and Combei et al. showed that ensemble approaches combining multiple SSL variants with different backend architectures can achieve superior results through late fusion techniques.

III. SAFE CHALLENGE OVERVIEW

The SAFE Challenge evaluates [33] audio deepfake detection systems across three distinct tasks designed to test different aspects of robustness.

- 1) **Task 1 (Generated Audio):** Detection of unmodified synthetic speech directly from TTS model output, testing basic spoofing detection capabilities.
- 2) **Task 2 (Processed Audio):** Detection of generated samples that have undergone compression and resampling operations, simulating real-world distribution scenarios.
- 3) **Task 3 (Laundered Audio):** Detection of deliberately processed synthetic audio designed to avoid detection systems, representing adversarial laundering attacks.

The SAFE evaluation dataset comprises human- and machine-generated speech audio tracks with several key characteristics that emphasize real-world applicability. Human speech samples are sourced from multiple origins and languages, ranging from high-quality studio recordings to lower-quality in-the-wild online recordings, as detailed in Table I under Task 1 - Real Audio sources. Machine-generated samples

²<https://github.com/imdatceleste/m-ailabs-dataset>

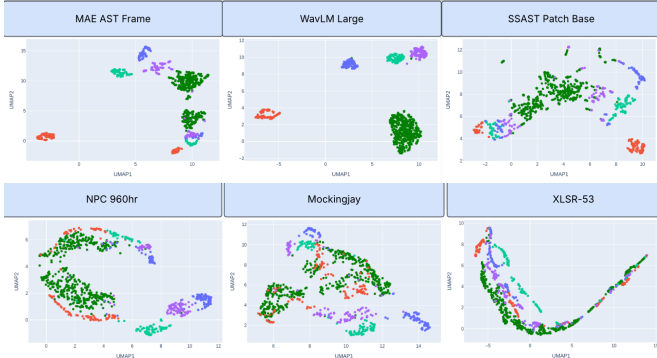


Fig. 1: UMAP visualizations of SSL model representations for ASVSpooF 2019 LA eval data. Top: Comparison of MAE AST Frame, WavLM Large, and SSAST Patch Base models. Bottom: Additional SSL model comparisons including NPC 960hr, Mockingjay, and XLSR-53. Green color represent bonafide, other colors represent various deepfakes.

are created using state-of-the-art TTS models, including both open-source and closed-source systems such as ElevenLabs, OpenAI, and various neural synthesis approaches (Table I, Task 1 - TTS Systems). Beyond unmodified audio detection, the challenge incorporates realistic degradation scenarios through Task 2, which applies various processing operations including compression codecs (AAC, MP3, Opus), resampling, pitch shifting, and additive noise to simulate real-world distribution conditions (Table I, Task 2 - Processed Audio). Task 3 addresses adversarial laundering attacks that deliberately process synthetic audio to evade detection systems, employing techniques such as playback through car environments, acoustic reverberation, and re-recording methods (Table I, Task 3 - Laundered Audio). Audio files vary in length up to 60 seconds with diverse compression formats, and the dataset maintains balance across different sources. Critically, the competition employs a fully blind evaluation protocol where no training data is released, ensuring that the systems must generalize from external training data to unknown test conditions.

IV. MODEL ARCHITECTURE AND SSL SELECTION

Our detection system employs a two-stage architecture combining self-supervised learning (SSL) front-ends with the AASIST (Audio Anti-Spoofing using Integrated Spectro-Temporal graph attention networks) back-end as proposed by [34]. To select optimal SSL models, we conducted UMAP visualization analysis on various SSL front-ends using ASVSpooF 2019 LA database. Figure 1 shows the UMAP visualization of MAE-AST Frame [35], WavLM Large [36], SSAST Patch Base [37], NPC 960hr [38], Mockingjay [39], and XLSR-53 models [40]. This analysis revealed that WavLM Large and MAE-AST Frame models provided the most discriminative feature representations for distinguishing between authentic and synthetic audio. The AASIST backend was chosen for its proven effectiveness in modeling both spectral and temporal spoofing artifacts.

TABLE I: SAFE Challenge Audio Sources by Task

Task	Category	Source Names
Task 1 (Unmodified)	Real Audio	Mandarin Podcast, FLEURS German, VSP Semi-professional, YouTube phonecall, VSP Documentary, Arabic Speech Corpus, High Quality Podcasts, Japanese Shortwave, Conference, English Podcast, FLEURS English, Djeco, Digitized Cassette, Librivox, Old Radio, phone home, Russian Audiobook, VSP Home Mic, Radio Drama, VSP Professional
	TTS Systems	elevenlabs, fish, hierspeech, kokoro, parler, seamless, style, cartesia, edge, f5, metavoices, openai, zonos
Task 2	Processed Audio	aac 16k, encodec, focalcodecs, mp3-aac-mp3 16k, mp3-aac 16k, mp3 16k, mp3 VBR, noise, opus 16k, phone audio, pitch shift, resample down/up, semanti-codec, snac, speech filter, time stretch, vorbis 16k
Task 3	Laundered Audio	car, played, played reverb car, reverb

V. DATASET INTEGRATION STRATEGY

A. TTS System Coverage Analysis

To systematically identify complementary datasets for optimal generalization performance, we conducted a comprehensive analysis of TTS system coverage across existing deepfake detection datasets. This analysis guides our strategic dataset selection approach. Table II provides a brief coverage of the TTS systems available in datasets used by this study.

1) *Complementarity Analysis*: Our analysis revealed three key insights: (1) **Technology Coverage Gaps** - most datasets focus on specific paradigms, with Famous Figures uniquely covering 2024-2025 TTS systems absent elsewhere; (2) **Optimal Combinations** - highest complementarity achieved by combining codec-based (CodecFake), multilingual (MLAAD), real-world noisy (SpoofCeleb), and cutting-edge (Famous Figures) approaches; (3) **Acoustic Condition Diversity** - combining clean studio, multilingual audiobook, and real-world noisy data provides comprehensive acoustic coverage.

2) *Strategic Implications*: This analysis demonstrates that systematic dataset combination based on complementary TTS coverage can address generalization challenges more effectively than single-dataset approaches, providing the foundation for our iterative experimental design.

B. Iterative Experimental Design

We systematically explored dataset integration through four iterative experiments, each building on the insights from the previous iterations. Our approach was motivated by the hypothesis that diverse multilingual datasets are crucial for generalization to unknown spoofing attacks.

1) *Iteration 1 (Baseline)*: We established a baseline using only the ASVSpooF 2019 LA training dataset, which comprises 25,380 samples (2,580 real, 22,800 fake). This single-dataset approach represents the conventional training paradigm and provided a reference point for measuring the impact of dataset diversification.

TABLE II: TTS Systems Coverage Across Deepfake Detection Datasets

Dataset	Total TTS Systems	Key TTS Architectures	Complementary Features
CodecFake	15 codec models (6 frameworks)	SpeechTokenizer, AcademicCodec, AudioDec, Encodec	Neural audio codec focus
ASVspoof 2019 LA	19 systems (6 known + 13 unknown)	WaveNet, Tacotron2, VAE-based, GMM-UBM	Traditional + neural mix
MLAAD	54-91 systems (21-42 architectures)	VITS variants, SpeechT5, multilingual models	Multilingual diversity
SpoofCeleb	23 systems	Contemporary TTS on VoxCeleb1	Real-world noisy conditions
Famous Figures	10 systems	StyleTTS2, XTTSv2, F5TTS, E2TTS, FishSpeech	Real-world, Latest 2024-2025 TTS

2) *Iteration 2 (Multi-dataset Integration)*: We expanded training data to include five complementary datasets with strategic sampling from each. For all datasets except ASVspoof 2019 LA, we applied an 80-20 split for training and validation data allocation.

- **ASVspoof 2019 LA**: Complete dataset with 25,380 samples providing established spoofing detection benchmarks
- **M-AILABS**: We sampled 20,000 multilingual authentic speech samples spanning 8 languages (English, French, German, Italian, Polish, Russian, Spanish, and Ukrainian). Following the 80-20 split, these samples were divided into 16,000 for training and 4,000 for validation, providing authentic multilingual speech baselines.
- **MLAAD (Multi-Language Audio Anti-Spoofing Dataset)**: We sampled 47,200 multilingual synthetic samples focusing specifically on languages available in M-AILABS plus Hindi language. This strategic language selection ensured consistency with our authentic multilingual data while covering diverse TTS architectures across multiple linguistic contexts.
- **CodecFake A2**: Since we were only interested in TTS systems, we specifically included only the A2 subset which uses VALL-E X, a modern neural codec-based TTS approach. We utilized 7,109 samples for training and 1,778 samples for validation, representing contemporary codec-based synthesis methods.
- We sampled 12,734 custom samples with 7,200 authentic and 5,534 synthetic samples from our Famous Figures dataset [23]. This dataset provides latest TTS systems for specific high-profile speakers.

This resulted in 108,423 training samples (25,780 authentic, 82,643 synthetic) and 45,606 validation samples, significantly expanding both linguistic and technical diversity compared to the single-dataset baseline.

3) *Iteration 3 (Audio Length Optimization)*: Using the same multi-dataset composition from Iteration 2, we increased audio segment length from 4 seconds to 12 seconds based on the observation that SAFE audio files can extend up to 60 seconds. This change was motivated by the hypothesis that longer temporal context would improve detection of complex spoofing artifacts that may require extended analysis windows.

C. Iteration 4 (Strategic Integration)

Based on our TTS system coverage analysis from Section V-A, which identified optimal complementarity combinations

for maximum technological and acoustic diversity, we refined the dataset composition through strategic optimization. This iteration specifically implements the findings that combining codec-based synthesis, multilingual diversity, real-world noisy conditions, and cutting-edge TTS technology provides the most comprehensive training foundation. The focus of Iteration 4 includes achieving balanced real/fake distribution, incorporating unseen languages and TTS systems, and ensuring comprehensive coverage across different synthesis paradigms.

SpoofCeleb Integration: Added 100,000 training samples (50,000 authentic, 50,000 synthetic) and 20,000 validation samples from this real-world noisy dataset, providing exposure to practical audio conditions that bridge the gap between clean laboratory data and real-world deployment scenarios.

MLAAD Language Selection: Strategically sampled 60,000 samples focusing on languages present in M-AILABS (English, German, Spanish, French, Italian, Polish, Russian, Ukrainian) plus Hindi. We implemented a careful train/validation split strategy to ensure both known and unknown TTS system coverage within each language:

- **English**: From 36 TTS systems, we randomly selected 7 systems for validation (2,000 samples) and used the remaining 29 systems for training (7,000 samples).
- **Ukrainian**: Used entirely for validation (5,000 samples), providing a completely unknown language condition for language generalization.
- **German, Spanish, French, and Italian**: For each language, we randomly selected 2 TTS systems for validation (2,000 samples each) with remaining systems allocated to training (7,000 samples for German, 6,000 samples each for Spanish, French, and Italian).
- **Polish**: Used 1 TTS system for validation (1000 samples) with remaining systems for training (5000 samples).
- **Russian**: Used entirely for training (5,000 samples), providing additional training diversity.
- **Hindi**: Dedicated entirely to training (2,000 samples), expanding linguistic coverage beyond European languages.

This careful partitioning resulted in 44,000 training samples and 16,000 validation samples, with validation splits containing both known architectures from unknown speakers and completely unknown language conditions (Ukrainian).

Famous Figures Refinement: Focused on 16,000 samples with balanced authentic/synthetic distribution, emphasizing Donald Trump and JD Vance for both categories while using other speakers only for authentic samples. This dataset

uniquely covers 2024-2025 TTS systems, which are absent from other datasets.

CodecFake Integration: Maintained 6,500 training and 1,600 validation samples from VALL-E X system to preserve codec-based synthesis representation.

M-AILABS Expansion: Increased to 60,000 samples (44,000 training, 16,000 validation) to better balance authentic multilingual content with the expanded synthetic data.

The final dataset comprised 200,000 training samples (101,200 authentic, 99,600 synthetic) and 56,600 validation samples (29,200 authentic, 27,400 synthetic), representing the most comprehensive multilingual, multi-domain training configuration guided by our systematic complementarity analysis. This configuration successfully achieves the targeted balanced real/fake distribution (approximately 50-50 split), incorporates multiple unseen languages and unknown TTS systems in validation splits, and provides comprehensive coverage across traditional neural TTS, codec-based synthesis, multilingual approaches, and cutting-edge 2024-2025 technologies.

VI. TRAINING DETAILS

All audio samples were resampled to 16 kHz and padded or cropped to fixed lengths depending on the iteration: 4 seconds for Iterations 1-2 and 12 seconds for Iterations 3-4. The longer audio segments in later iterations were chosen to provide extended temporal context for detecting complex spoofing artifacts, given that SAFE challenge audio files can extend up to 60 seconds.

We applied RawBoost data augmentation [29] across all iterations, using linear and non-linear convolutive noise combined with impulsive signal-dependent additive noise strategies optimal for logical access scenarios. All models were trained using Adam optimizer with learning rate of 10^{-6} , and binary cross-entropy loss with class weighting. We used S3PRL toolkit [41], [42] to extract SSL embeddings: WavLM (1024 dimensions) and MAE-AST Frame (768 dimensions), both fed to 128-dimensional fully connected layers before AASIST back-end classification. Training was conducted separately for each iteration over 50 epochs on Nvidia A100 GPU, with reproducible results available through open source code³.

VII. RESULTS AND ANALYSIS

Table III presents the performance progression across our four iterative experiments on both the SAFE Challenge tasks and the ITW benchmark. The results demonstrate the systematic impact of dataset integration strategies on detection performance across different evaluation scenarios.

Iteration 1 (Baseline): Using traditional AASIST trained solely on ASVspoof 2019 LA with 4-second audio segments, we achieved baseline performance of 0.531, 0.589, and 0.492 balanced accuracy on Tasks 1, 2, and 3 respectively.

Iteration 2 (Multi-dataset Integration): Incorporating five complementary datasets (ASVspoof 2019 LA, M-AILABS, MLAAD, CodecFake A2, and Famous Figures) with SSL

TABLE III: Performance Progression Across Four Iterations on SAFE Challenge Tasks and ITW Benchmark.

Iter	SSL Model	SAFE Challenge (BA)			ITW Benchmark	
		Task 1	Task 2	Task 3	BA	EER (%)
1	AASIST	0.531	0.589	0.492	0.616	35.61
2	WavLM	0.745	0.587	0.478	0.875	8.46
2	MAE-AST	0.607	0.587	0.597	0.648	24.79
3	WavLM	0.766	0.765	0.518	0.856	12.05
4	WavLM	0.810	0.819	0.496	0.905	8.42
4	MAE-AST	0.640	0.536	0.623	0.603	39.9

front-ends yielded substantial improvements. WavLM Large achieved 0.745, 0.587, and 0.478 on Tasks 1, 2, and 3, representing a 40.3% improvement on Task 1. MAE-AST Frame showed 0.607, 0.587, and 0.597, demonstrating competitive performance on Task 3.

Iteration 3 (Audio Length Optimization): Extending audio segments from 4 to 12 seconds using the same multi-dataset composition produced notable gains for WavLM Large: 0.766 (Task 1), 0.765 (Task 2), and 0.518 (Task 3). The 30.3% improvement in Task 2 performance specifically validates the importance of longer temporal context for detecting processed audio artifacts.

Iteration 4 (Strategic Integration): Our final configuration, incorporating SpoofCeleb, refined dataset balancing and strategic train-validation splits, achieved optimal performance with WavLM Large: 0.810 (Task 1), 0.819 (Task 2), and 0.496 (Task 3). This represents cumulative improvements of 52.5% and 39.0% for Tasks 1 and 2 respectively from the baseline. Task 3 performance remained stable around 0.49-0.52 across iterations, reflecting the inherent challenge of detecting adversarially laundered audio.

A. SAFE Challenge Performance

Our best configuration (Iteration 4, WavLM Large) demonstrated strong performance across both evaluation phases of the SAFE Challenge. On the public leaderboard, we achieved second place across all three tasks (Tasks 1, 2, and 3). On the private leaderboard, which is a superset of the public split, we secured third place for all tasks. This consistent top-tier ranking among international research teams validates the effectiveness of our multilingual dataset integration approach. The SAFE Challenge evaluation strategy involved selecting the best-performing model architecture for each task based on the submitted model. WavLM Large model (iteration 4) proved optimal for Tasks 1 and 2. For Task 3 (laundered audio), MAE-AST Frame’s more robust temporal modeling against adversarial processing made it the preferred choice, demonstrating the value of architectural diversity in handling different threat scenarios.

B. Task-Specific Analysis

Our evaluation reveals distinct challenges across the three SAFE Challenge tasks. Task 1 (unmodified audio) shows fundamental detection challenges with WavLM Large achieving strong performance (BA = 0.810) while MAE-AST Frame reaches 0.640. Task 2 (processed audio) demonstrates our

³<https://github.com/issflab/ssl-antispoofing/>

highest performance levels (BA = 0.819 with WavLM Large), showing resilience to compression and resampling operations. Task 3 (laundered audio) presents the greatest challenge, where WavLM Large shows consistent degradation (0.478-0.518 range) compared to Tasks 1-2, while MAE-AST Frame maintains more stable performance across all tasks (0.597-0.640 range for iteration 2, 0.536-0.640 for iteration 4).

TABLE IV: Balanced Accuracy on Task 1 for Pristine and Generated Sources (Iteration 4, WavLM Large). Bold values indicate balanced accuracy less than or equal to 0.60.

Pristine 1	Acc.	Pristine 2	Acc.	Generated	Acc.
MandPod1	0.87	EngPod	0.86	elevenlabs	0.64
FleurGer	0.87	FleurEng	0.86	fish	0.81
VSPSemi	0.58	DigCass	0.71	hierspeech	0.87
YTPPhone	0.87	Dipco	0.84	kokoro	0.87
VSPDoc	0.84	Librivox	0.69	parler	0.86
ArabCorpus	0.87	OldRadio	0.69	seamless	0.76
HQPod	0.83	PhoneHome	0.53	style	0.86
JapSWave	0.39	RussAudiobook	0.52	cartesia	0.47
Conf	0.48	MandPod2	0.86	f5	0.63
VSPHomeMic	0.87	RadioDrama	0.81	metavox	0.58
VSPProf	0.84			zonos	0.65

Legend: MandPod1/2 = Mandarin Podcast 1/2; FleurGer/Eng = Fleurs German/English; VSPSemi/Prof = VSP Semi-professional/Professional; YTPPhone = YouTube phonecall; VSPDoc = VSP Documentary; HQPod = High Quality Podcasts; JapSWave = Japanese Shortwave

TABLE V: Balanced Accuracy for Task 2 (Processed) and Task 3 (Laundered) Sources (Iteration 4, WavLM Large). Bold values indicate balanced accuracy less than or equal to 0.60.

Processed 1		Processed 2		Laundered	
Source	Task 2	Source	Task 2	Source	Task 3
aac 16k	0.80	pitch shift	0.81	car	0.51
encodec	0.84	resample ↓	0.77	played	0.67
focallcodec	0.83	resample ↑	0.76	reverb	0.64
mp3-aac-mp3	0.84	sem-codec	0.74	all 3	0.49
mp3-aac 16k	0.81	snac	0.73		
mp3 16k	0.79	speech filt.	0.79		
mp3 VBR	0.79	time stret.	0.85		
noise	0.52	vorbis 16k	0.75		
opus 16k	0.73	phone audio	0.85		

C. Source-Level Performance Analysis

Source-level analysis was conducted for each major SAFE Challenge task using the Iteration 4 - WavLM Large system. The results for Task 1 (unmodified audio) are summarized in Table IV, which lists balanced accuracy for both pristine and generated sources. Tasks 2 and 3 (processed audio and laundered audio, respectively) are reported in Table V. In these tables, bold values indicate cases where balanced accuracy was ≤ 0.60 , highlighting challenging sources for our system.

As seen in Table IV, most pristine (real) sources exhibit high detection accuracy (≥ 0.80), notably *MandPod1*, *FleurGer*, *EngPod*, and *VSPHomeMic*. However, several real sources stand out as considerably more challenging, with balanced accuracy of ≤ 0.60 : *VSPSemi* (**0.58**), *JapSWave* (**0.39**), *Conf* (**0.48**), *PhoneHome* (**0.53**), and *RussAudiobook* (**0.52**).

This variability points to vulnerabilities in detecting authentic, diverse, or lower-quality real-world audio.

The generated sources (right-most column of Table IV) generally achieve high detection rates as well. However, certain synthetic systems are more difficult to detect, including *cartesia* (**0.47**) and *metavox* (**0.58**). Most mainstream approaches, such as *hierspeech*, *kokoro*, and *openai*, are robustly detected (> 0.80), while the worst-performing synthetic and pristine sources may warrant special attention for system improvement.

Table V shows the balanced accuracy for each processed and laundered source. For processed sources (Task 2), the majority maintain balanced accuracy above 0.70, with a few notable exceptions: the **noise** condition (**0.52**) significantly lowers accuracy, showing that added noise can effectively degrade detection reliability. All other processed sources—various codecs, resampling, pitch shift, semantic encoding—sustain robust detection (typically 0.73 – 0.85).

Task 3 (laundered audio) remains challenging, with the lowest performance: **car** (**0.51**) and **played reverb car** (**0.49**) both fall below the 0.60 threshold, and even the best-performing laundered source (*played*, 0.67) shows considerable performance degradation compared to unprocessed tasks.

D. ITW Evaluation

Performance on the In-The-Wild (ITW) benchmark provides crucial validation of our approach’s generalization capabilities to real-world deepfakes from social media platforms. Our results demonstrate substantial improvements across iterations, with balanced accuracy progressing from 0.616 (baseline) to 0.905 (Iteration 4, WavLM Large) and EER reducing from 35.61% to 8.42%. Each iteration shows significant improvement: multi-dataset integration (BA: 0.616 \rightarrow 0.875), audio length optimization (BA: 0.856), and strategic integration achieving optimal performance (BA: 0.905, EER: 8.42%).

VIII. CONCLUSION

This work demonstrates that strategic integration of multilingual data sets significantly improves detection of deepfake audio in diverse scenarios. Our systematic approach, combining six complementary datasets with SSL-based architectures, achieved competitive SAFE Challenge rankings (2nd-3rd place) and substantial ITW benchmark improvements (EER: 35.61% \rightarrow 8.42%).

Our key findings include: (1) dataset diversity provides dramatic performance gains over single-dataset approaches, (2) WavLM Large and MAE-AST Frame offer complementary strengths for different threat scenarios, (3) longer audio segments improve detection of processed artifacts, and (4) laundered audio presents critical challenges with false positive vulnerabilities that require careful deployment consideration.

Our source-level analysis reveals specific failure patterns and architectural trade-offs that inform future research directions. Although substantial progress has been made on clean and processed audio detection, adversarial laundering remains a significant challenge requiring continued investigation.

REFERENCES

- [1] A. Nautsch, X. Wang, N. Evans, T. H. Kinnunen, V. Vestman, M. Todisco, H. Delgado, M. Sahidullah, J. Yamagishi, and K. A. Lee, "Asvspoof 2019: spoofing countermeasures for the detection of synthesized, converted and replayed speech," *IEEE Transactions on Biometrics, Behavior, and Identity Science*, vol. 3, no. 2, pp. 252–265, 2021.
- [2] X. Wang, J. Yamagishi, M. Todisco, H. Delgado, A. Nautsch, N. Evans, M. Sahidullah, V. Vestman, T. Kinnunen, K. A. Lee *et al.*, "Asvspoof 2019: A large-scale public database of synthesized, converted and replayed speech," *Computer Speech & Language*, vol. 64, p. 101114, 2020.
- [3] M. R. Kamble, H. B. Sailor, H. A. Patil, and H. Li, "Advances in anti-spoofing: from the perspective of asvspoof challenges," *APSIPA Trans. on Signal and Information Processing*, vol. 9, p. e2, 2020.
- [4] N. M. Müller, P. Czempin, F. Dieckmann, A. Froghyar, and K. Böttinger, "Does Audio Deepfake Detection Generalize?" Apr. 2022, arXiv:2203.16263 [cs, eess]. [Online]. Available: <http://arxiv.org/abs/2203.16263>
- [5] X. Liu, X. Wang, M. Sahidullah, J. Patino, H. Delgado, T. Kinnunen, M. Todisco, J. Yamagishi, N. Evans, A. Nautsch, and K. A. Lee, "ASVspoof 2021: Towards Spoofed and Deepfake Speech Detection in the Wild," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 31, pp. 2507–2522, 2023. [Online]. Available: <https://ieeexplore.ieee.org/document/10155166/>
- [6] J. Yamagishi, X. Wang, M. Todisco, M. Sahidullah, J. Patino, A. Nautsch, X. Liu, K. A. Lee, T. Kinnunen, N. Evans *et al.*, "Asvspoof 2021: accelerating progress in spoofed and deepfake speech detection," in *ASVspoof 2021 Workshop-Automatic Speaker Verification and Spoofing Countermeasures Challenge*, 2021.
- [7] X. Wang, H. Delgado, H. Tak, J. weon Jung, H. jin Shim, M. Todisco, I. Kukanov, X. Liu, M. Sahidullah, T. H. Kinnunen, N. Evans, K. A. Lee, and J. Yamagishi, "Asvspoof 5: crowdsourced speech data, deepfakes, and adversarial attacks at scale," in *The Automatic Speaker Verification Spoofing Countermeasures Workshop (ASVspoof 2024)*, 2024, pp. 1–8.
- [8] Z. Wu, T. Kinnunen, N. Evans, and J. Yamagishi, "Asvspoof 2015: Automatic speaker verification spoofing and countermeasures challenge evaluation plan," *Training*, vol. 10, no. 15, p. 3750, 2014.
- [9] 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," *International Joint Conference on Neural Networks (IJCNN)*, 2024.
- [10] J.-w. Jung, Y. Wu, X. Wang, J.-H. Kim, S. Maiti, Y. Matsunaga, H.-j. Shim, J. Tian, N. Evans, J. S. Chung, W. Zhang, S. Um, S. Takamichi, and S. Watanabe, "SpoofCeleb: Speech Deepfake Detection and SASV In The Wild," Sep. 2024, arXiv:2409.17285 [cs]. [Online]. Available: <http://arxiv.org/abs/2409.17285>
- [11] L. Zhang, H. Zhao, Q. Meng, Y. Chen, M. Liu, and L. Xie, "Beijing ZKJ-NPU Speaker Verification System for VoxCeleb Speaker Recognition Challenge 2021," arXiv:2109.03568 [cs, eess], Nov. 2021, arXiv: 2109.03568. [Online]. Available: <http://arxiv.org/abs/2109.03568>
- [12] H. Ali, S. Subramani, S. Sudhir, R. Varahamurthy, and H. Malik, "Is audio spoof detection robust to laundering attacks?" in *Proceedings of the 2024 ACM Workshop on Information Hiding and Multimedia Security*, 2024, pp. 283–288.
- [13] Y. Xie, Y. Lu, R. Fu, Z. Wen, Z. Wang, J. Tao, X. Qi, X. Wang, Y. Liu, H. Cheng *et al.*, "The codefake dataset and countermeasures for the universally detection of deepfake audio," arXiv preprint arXiv:2405.04880, 2024.
- [14] X. Zhang, D. Zhang, S. Li, Y. Zhou, and X. Qiu, "Speechtokenizer: Unified speech tokenizer for speech language models," in *The Twelfth International Conference on Learning Representations*, 2024.
- [15] A. Défossez, J. Copet, G. Synnaeve, and Y. Adi, "High fidelity neural audio compression," arXiv preprint arXiv:2210.13438, 2022.
- [16] Z. Du, S. Zhang, K. Hu, and S. Zheng, "Funcodec: A fundamental, reproducible and integrable open-source toolkit for neural speech codec," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 591–595.
- [17] 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)*. IEEE, 2024, pp. 921–928.
- [18] V. Popov, I. Vovk, V. Gogoryan, T. Sadekova, and M. Kudinov, "Grad-tts: A diffusion probabilistic model for text-to-speech," in *International conference on machine learning*. PMLR, 2021, pp. 8599–8608.
- [19] K. Shen, Z. Ju, X. Tan, Y. Liu, Y. Leng, L. He, T. Qin, S. Zhao, and J. Bian, "Naturalspeech 2: Latent diffusion models are natural and zero-shot speech and singing synthesizers," arXiv preprint arXiv:2304.09116, 2023.
- [20] Y. A. Li, C. Han, V. Raghavan, G. Mischler, and N. Mesgarani, "Styletts 2: Towards human-level text-to-speech through style diffusion and adversarial training with large speech language models," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [21] S. Mehta, R. Tu, J. Beskow, É. Székely, and G. E. Henter, "Matcha-tts: A fast tts architecture with conditional flow matching," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 11 341–11 345.
- [22] S. Kim, K. Shih, J. F. Santos, E. Bakhturina, M. Desta, R. Valle, S. Yoon, B. Catanzaro *et al.*, "P-flow: a fast and data-efficient zero-shot tts through speech prompting," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [23] H. Ali, S. Subramani, R. Varahamurthy, N. Adupa, L. Bollinani, and H. Malik, "Collecting, curating, and annotating good quality speech deepfake dataset for famous figures: Process and challenges," arXiv preprint arXiv:2507.00324, 2025.
- [24] V. Elliott, "The Biden Deepfake Robocall Is Only the Beginning," *Wired*, Jan. 2024, section: tags. [Online]. Available: <https://www.wired.com/story/biden-robocall-deepfake-danger/>
- [25] M. Spring, "Sadiq Khan says fake AI audio of him nearly led to serious disorder," Feb. 2024. [Online]. Available: <https://www.bbc.com/news/uk-68146053>
- [26] E. Casanova, K. Davis, E. Gölge, G. Gökmar, I. Gulea, L. Hart, A. Aljafari, J. Meyer, R. Morais, S. Olayemi *et al.*, "Xtts: a massively multilingual zero-shot text-to-speech model," arXiv preprint arXiv:2406.04904, 2024.
- [27] Y. Chen, Z. Niu, Z. Ma, K. Deng, C. Wang, J. Zhao, K. Yu, and X. Chen, "F5-tts: A fairytale that fakes fluent and faithful speech with flow matching," arXiv preprint arXiv:2410.06885, 2024.
- [28] S. E. Eskimez, X. Wang, M. Thakker, C. Li, C.-H. Tsai, Z. Xiao, H. Yang, Z. Zhu, M. Tang, X. Tan *et al.*, "E2 tts: Embarrassingly easy fully non-autoregressive zero-shot tts," in *2024 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2024, pp. 682–689.
- [29] H. Tak, M. Kamble, J. Patino, M. Todisco, and N. Evans, "Rawboost: A raw data boosting and augmentation method applied to automatic speaker verification anti-spoofing," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 6382–6386.
- [30] O. Pascu, A. Stan, D. Oneata, E. Oneata, and H. Cucu, "Towards generalisable and calibrated audio deepfake detection with self-supervised representations," in *Interspeech*, vol. 2024, 2024, pp. 4828–4832.
- [31] T. Stourbe, V. Miara, T. Lepage, and R. Dehak, "Exploring wavlm backends for speech spoofing and deepfake detection," in *Proc. ASVspoof 2024*, 2024, pp. 72–78.
- [32] D. Combei, A. Stan, D. Oneata, and H. Cucu, "Wavlm model ensemble for audio deepfake detection," arXiv preprint arXiv:2408.07414, 2024.
- [33] T. Kirill, P. Cummer, P. Pherwani, J. Aslam, M. Davinroy, P. Bautista, L. Cassani, and M. Stamm, "Safe: Synthetic audio forensics evaluation challenge," in *Proceedings of the ACM Workshop on Information Hiding and Multimedia Security*, 2025, pp. 174–180.
- [34] H. Tak, M. Todisco, X. Wang, J.-w. Jung, J. Yamagishi, and N. Evans, "Automatic speaker verification spoofing and deepfake detection using wav2vec 2.0 and data augmentation," in *The Speaker and Language Recognition Workshop*, 2022.
- [35] A. Baade, P. Peng, and D. Harwath, "Mae-ast: Masked autoencoding audio spectrogram transformer," arXiv preprint arXiv:2203.16691, 2022.
- [36] 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, and F. Wei, "Wavlm: Large-scale self-supervised pre-training for full stack speech processing," 2021.
- [37] Y. Gong, C.-I. Lai, Y.-A. Chung, and J. Glass, "Ssast: Self-supervised audio spectrogram transformer," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 10, 2022, pp. 10 699–10 709.
- [38] A. H. Liu, Y.-A. Chung, and J. Glass, "Non-autoregressive predictive coding for learning speech representations from local dependencies," arXiv preprint arXiv:2011.00406, 2020.

- [39] D. B. T. Encoders, “Mockingjay: Unsupervised speech representation learning with,” 2020.
- [40] A. Babu, C. Wang, A. Tjandra, K. Lakhotia, Q. Xu, N. Goyal, K. Singh, P. von Platen, Y. Saraf, J. Pino *et al.*, “Xls-r: Self-supervised cross-lingual speech representation learning at scale,” *arXiv preprint arXiv:2111.09296*, 2021.
- [41] S. wen Yang, P.-H. Chi, Y.-S. Chuang, C.-I. J. Lai, K. Lakhotia, Y. Y. Lin, A. T. Liu, J. Shi, X. Chang, G.-T. Lin, T.-H. Huang, W.-C. Tseng, K. tik Lee, D.-R. Liu, Z. Huang, S. Dong, S.-W. Li, S. Watanabe, A. Mohamed, and H. yi Lee, “SUPERB: Speech Processing Universal PERformance Benchmark,” in *Proc. Interspeech 2021*, 2021, pp. 1194–1198.
- [42] S.-w. Yang, H.-J. Chang, Z. Huang, A. T. Liu, C.-I. Lai, H. Wu, J. Shi, X. Chang, H.-S. Tsai, W.-C. Huang *et al.*, “A large-scale evaluation of speech foundation models,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.