


Multi-source Domain Adaptation for Text-independent Forensic Speaker Recognition

Zhenyu Wang, and John H. L. Hansen , *Fellow, IEEE*

Abstract—Adapting speaker recognition systems to new environments is a widely-used technique to improve a well-performing model learned from large-scale data towards a task-specific small-scale data scenarios. However, previous studies focus on single domain adaptation, which neglects a more practical scenario where training data are collected from multiple acoustic domains needed in forensic scenarios. Audio analysis for forensic speaker recognition offers unique challenges in model training with multi-domain training data due to location/scenario uncertainty and diversity mismatch between reference and naturalistic field recordings. It is also difficult to directly employ small-scale domain-specific data to train complex neural network architectures due to domain mismatch and performance loss. Fine-tuning is a commonly-used method for adaptation in order to retrain the model with weights initialized from a well-trained model. Alternatively, in this study, three novel adaptation methods based on domain adversarial training, discrepancy minimization, and moment-matching approaches are proposed to further promote adaptation performance across multiple acoustic domains. A comprehensive set of experiments are conducted to demonstrate that: 1) diverse acoustic environments do impact speaker recognition performance, which could advance research in audio forensics, 2) domain adversarial training learns the discriminative features which are also invariant to shifts between domains, 3) discrepancy-minimizing adaptation achieves effective performance simultaneously across multiple acoustic domains, and 4) moment-matching adaptation along with dynamic distribution alignment also significantly promotes speaker recognition performance on each domain, especially for the LENA-field domain with noise compared to all other systems. Advancements shown here in adaptation therefore help ensure more consistent performance for field operational data in audio forensics.

Index Terms—discrepancy loss, forensics, multi-source domain adaptation, domain adversarial training, maximum mean discrepancy, moment-matching, speaker recognition.

I. INTRODUCTION

IN general, no two speakers are identical, differing in anatomy, physiology, and acoustically from a speech production viewpoint. Considering human speech as a discriminative biometric, speaker recognition serves as an important tool in law enforcement, national security, and forensics in general. The need for forensic speaker recognition arises

when an individual contributes his/her voice as evidence, including telephone recordings, wiretaps, audio surveillance, or informant recordings [1]. The use of technology for forensic speaker recognition has been considered as early as 1926 based on speech waveform analysis [2]. It was popularized much later in the 1970s, when it came to be known incorrectly as the “voiceprint” [3]. Approaches to forensic speaker recognition include spectrographic, auditory, acoustic-phonetic, and automatic. Forensic speaker recognition is commonly performed fully or partially by human expert phoneticians who generally have backgrounds in linguistics and statistics. Full or assisted automatic approaches are also considered as an efficient tool for forensic speaker recognition to aid the forensic examiner in quantifying the strength of evidence [4]–[7].

In the forensic context, speaker recognition assists investigators and legal/courtrooms (judge or jury) to identify an unknown speaker suspected of a crime in legal proceedings. In general, for forensic speaker recognition, a likelihood ratio is needed to determine how likely a voice recording was produced by a speaker of known identity (typically a suspect) or not [2], [8].

Great progress has been made in speaker recognition in recent decades, thus solidifying automatic speaker recognition as a core tool in the forensic field. Previously, the segment-level vectors that represent speech entitled i-Vectors with probabilistic linear discriminant analysis (PLDA) as a backend have dominated the text-independent speaker recognition research field [9]. Additionally, i-Vector variants have been widely used in multiple fields of paralinguistic speech attribute recognition [10]–[12]. With the emergence of large speaker labeled audio datasets and growing computational resources, there is increased interest in applying more effective approaches including x-Vector and other neural network architectures to speaker recognition tasks [13]–[19].

Forensic speech data as potential evidence can be obtained in random naturalistic environments resulting in variable data quality. Additionally, speaker-based uncertainties such as stress, sentiment, vocal effort, and other intrinsic speaker factors introduce unknown mismatch challenges [1]. Mismatch variability consisting of intrinsic and extrinsic characteristics can degrade performance of speaker recognition. Intrinsic speaker characteristics represent speech traits that are dependent on the speaker vs. extrinsic characteristics that are dependent on audio capture and environmental factors [1]. Intrinsic properties include the speaker’s age, gender, ethnicity, vocal effort, noise-induced Lombard effect [20], situational stress, emotional and physical state (e.g., angry, sad, stressed, distracted, etc.). Extrinsic properties include

Manuscript received April 1, 2021; revised July 27, 2021 and November 8, 2021; accepted November 15, 2021. Date of publication November 26, 2021; date of current version December 20, 2021. This work was supported by the University of Texas at Dallas from the Distinguished University Chair in Telecommunications Engineering held by John H. L. Hansen. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Alberto Abad. (Corresponding author: John H. L. Hansen.)

The authors are with Center for Robust Speech Systems, Erik Jonsson School of Engineering University of Texas at Dallas, Richardson, TX 75080, USA (e-mail: zhenyu.wang@utdallas.edu; john.hansen@utdallas.edu).

Digital Object Identifier 10.1109/TASLP.2021.3130975

recording equipment conditions such as microphone type and placement, background noise, room reverberation, and other environmental scenario-based issues. Some factors, such as noise and non-target speech, may impact system performance by their mere presence. Variations caused by multi-faceted acoustic scenarios pose major challenges for effective model development to recognize a speaker in mismatched conditions.

To mitigate the impact of condition mismatch or domain shift, domain adaptation is needed to generalize a well-trained model learned for one acoustic domain to encompass new domains with task-specific data. Some theoretical insights can be drawn from the computer vision (CV) field. For example, some domain adaptation methods [21]–[23] employ neural networks with a Maximum Mean Discrepancy (MMD) loss to diminish domain discrepancy; other methods introduce novel model training schema to optimize domain alignments such as moment matching [24], adversarial domain confusion [25]–[27] and Generative Adversarial Network (GAN)-based alignment [28]–[30]. Such methods have also been adopted by the speaker adaptation community to address tasks focusing on domain mismatch, training techniques, and new architectures. Learning domain-invariant speaker embeddings with GAN have been considered to mitigate the impact of language mismatch between training and evaluation sets [31], [32]. For example, [33], [34] proposed to adapt speaker recognition systems with MMD at both frame-level and segment-level, which is assumed to be more robust against duration discrepancy. New architectures optimizing with objective-targeting losses have also been investigated to extract robust embeddings [35], [36]. However, the domain adaptation methods noted above mainly align the distributions of representations that may only contain partial information for a single domain, which are not practical in many forensic scenarios where speech samples are typically collected from diverse naturalistic domains with multiple mismatches containing unknown context knowledge, thereby, requiring Multi-source domain adaptation (MSDA). Our previous work [37] utilized the discrepancy minimization method for cross-domain adaptation and evaluated system performance within a closed-set using forensic data. Related works have also been applied in the CV field recently, such as novel cross-domain structures based on the formalism of multi-class domain adaptation were proposed. These studies consider the concept of minimizing the domain distance or category shift with measures such as MMD, Moment Distance (MD), and Multi-Class Scoring Disagreement (MCSD) [38]–[42].

In this study, we develop three multi-source domain adaptation approaches to learn domain-invariant information across naturalistic environments containing extrinsic variations. These variations alter speaker identity traits, whose instantiation variants require new learning objectives which either coincide with or resemble widely-used methods, thus partially underscoring their effectiveness in more practical scenarios. To address the lack of available real naturalistic forensic audio corpora with ground-truth speaker identity, we introduce our CRSS-Forensic dataset for benchmarking state-of-the-art multi-source domain adaptation methods. The dataset includes four subsets: (i) Clean (e.g. audio recorded

with a close-talk mic and a desk-top mic), (ii) Far-field (e.g. audio recorded with distance mics), (iii) LENA-booth (e.g. audio recorded with a asynchronous mobile data collection platform called LENA worn by the participant), and (iv) LENA-field (e.g. audio recorded in public environments with noise), where three kinds of mismatch such as distance mismatch, channel mismatch, and noise mismatch exist among these subsets. First, an x-Vector system is pre-trained with a large-scale VoxCeleb dataset, followed by fine-tuning the high-level neural network layers to learn speaker information from the CRSS-Forensic corpus. In addition to the pre-trained x-Vector model, we perform a multi-source domain adaptation using two alternative methods. One is based on discrepancy minimization to align the domain-specific distributions with maximum mean discrepancy (MMD). The second employs a moment-matching method to minimize the inter-domain discrepancies and dynamically aligns the moments of embedding distributions with an adversarial training strategy. In terms of our test protocol, we evaluate the pre-trained x-Vector system, fine-tuned system, discrepancy-minimization adaptation system, and moment-matching system with the Phase-I portion of the CRSS-Forensic dataset under an open-set framework. The main contributions of this study are as follows,

- 1) we demonstrate the impact of different acoustic environments on speaker recognition system performance.
- 2) A set of speaker recognition adaptation approaches are proposed to address forensic speaker recognition under diverse acoustic environments.
- 3) A discussion regarding best practices on how the proposed speaker recognition system could assist the "trier of fact" (i.e., a judge, a panel of judges, or a jury) in making decisions regarding the origin of speech on voice recordings of a speaker whose identity is in question.

This paper is organized as follows: Sec. II describes the x-Vector backbone system and fine-tuning details. Sec. III elaborates on our domain adversarial training approach. The description of our proposed adaptation system framework based on discrepancy minimization along with each component are presented in Sec. IV. Sec. V elaborates on the proposed system framework based on moment matching and our corresponding adversarial training schema. A brief description of each system's evaluation corpus and configurations are also illustrated, and the dataset description is included in Sec. VI. The effectiveness of the proposed methods is demonstrated in Sec. VII using a performance comparison across each sub-domain subset of the CRSS-Forensic corpus. Finally, conclusions are summarized in Sec. VIII.

II. BACKBONE SYSTEM AND FINE-TUNING

Since the x-Vector [13] has shown competitive results when trained on large proprietary datasets and is widely accepted as an effective speaker recognition solution, we employ the x-Vector system as the backbone system. Here, we pre-train an x-Vector system with large-scale VoxCeleb data to obtain a preliminary discriminative speaker representation, then fine-

tune that model using audio data from the CRSS-Forensic corpus.

A. Pre-trained System

The time-delay layer $\mathcal{F}(N, D, K)$ [43] forms the basic component of the x-Vector system which computes fixed-length speaker embeddings from variable-length acoustic segments. At each time-step, activations from the previous layer are computed using a context width of K , and a dilation of D . Here, N represents the output embedding dimension. The temporal short-term context is processed by this feed-forward time-delay architecture at the frame-level, where the statistics pooling layer is used to aggregate over frame-level representations to compute corresponding mean and standard deviations as the concatenation output. At the segment-level, additional fully-connected layers are used to operate on non-temporal concatenated information followed by a softmax output layer [44]. Ultimately, the goal of this architecture is to generate speaker embeddings over the entire utterance that hopefully will generalize well to unseen speakers within the training set. Therefore, suppose there are S speakers with M training samples, then the training objective is to maximize the probability $P(y_S | \mathbf{x}_{1:T}^{(M)})$ for speaker S given the T input frames $\mathbf{x}_1^{(M)}, \mathbf{x}_2^{(M)}, \dots, \mathbf{x}_T^{(M)}$. The optimization process can then be written as,

$$\Theta^* = \arg \max_{\Theta} [\mathbb{E}_{\mathbf{x}} [\log p_{\Theta}(y | \mathbf{x})]]. \quad (1)$$

Here, $\Theta = \{\mathbf{W}^l, \mathbf{b}^l\}_{l=1}^L$ denotes the trainable parameters of the L -layer neural network, (\mathbf{x}, y) represents the frame-level feature of an utterance and its corresponding label.

In this study, we employ the same backend probabilistic linear discriminant analysis (PLDA) [45] training method for each system. Speaker embeddings are centered and dimensionality reduction is accomplished using Linear Discriminant Analysis (LDA). After LDA, scores for pairs of length-normalized embeddings are generated using the PLDA model and normalized with an adaptive s-norm [46]. Next, a PLDA backend is used to compute scores for paired embeddings, which enables a similarity metric to be trained on potentially diverse situational datasets.

B. Fine-tuning

In general, x-Vector performance appears to be highly sensitive to both the amount and type of training data. This deep neural network has an extensive number of parameters, which for our solution totals 4.4 million parameters excluding the softmax output layer (it is not needed after training and will of course vary across different tasks). Any forensic dataset to be examined, in general, might not be very large in size due to the specialty of how the forensic dataset was collected in restricted acoustic environments. Additionally, speech samples may include variability due to vocal effort such as whisper-to-shout over 911 emergency calls, whereas others might include situational stress in a field location or interview room [1], [47]. Training the x-Vector model on a small or domain-mismatched dataset greatly affects the model's ability to generalize, often

resulting in over-fitting, especially if the last few layers of the network are fully connected layers. Therefore, model adaptation with fine-tuning is indispensable in this case. More often in practice, existing networks trained on a large dataset such as VoxCeleb [48] would continue to be trained on a targeted smaller task-specific dataset. Given that if the small dataset is not drastically different in terms of context from the original training dataset, the pre-trained model is assumed to have already learned basic speaker-based features relevant to a target final task. There is a common practice to truncate the last layer (softmax layer) of the pre-trained network and replace it with a new softmax layer consistent with speaker labels of a new task to adapt the network. Since we expect pre-trained weights to be quite effective compared to randomly initialized weights, it is important not to modify them either too quickly or too much. Thus, an initial fine-tuning learning rate should be smaller than the one used for training from scratch. Additionally, it is beneficial to freeze the weights of the first few layers of the pre-trained network. Since the first few layers capture universal acoustic features that are also relevant to those for a new task, instead of keeping weights intact, we encourage the network to focus on learning task-specific features for the subsequent intermediate and high-level layers.

III. DOMAIN ADVERSARIAL TRAINING

The commonly-used Fine-tuning method is generally effective in single-source domain adaptation, however, it cannot address the loss in performance resulting from a domain shift. Optimizing the classification objectives alone cannot guarantee effective generalization to multiple domains simultaneously without reducing the divergence between diverse distributions for each domain. Domain adversarial training utilizes domain information and promotes the emergence of features that are discriminative for speaker identity and invariant with respect to domain shift [26], [49].

We decompose a deep feed-forward architecture (see in Fig. 1) into three parts including (i) a universal feature extractor G with the parameters θ_g , (ii) a speaker label classifier C_s with the parameters θ_s , and a (iii) domain label classifier C_d with parameters θ_d . As training progresses, the parameters θ_g maximize the loss of the domain label classifier while simultaneously the parameters θ_d minimize the loss of the domain label classifier. Additionally, the parameters θ_s minimize the loss of the speaker label classifier. The multi-task softmax cross-entropy \mathcal{J} loss can be written as,

$$\begin{aligned} Loss_{DAT} = \sum_{i=1}^N \sum_{j=1}^{M^i} & \left(\mathcal{J}(C_s(G(\mathbf{x}_j^i)), y_j^i) \right. \\ & \left. - \lambda \mathcal{J}(C_d(G(\mathbf{x}_j^i)), d_j^i) \right). \end{aligned} \quad (2)$$

Given N domains, here $(\mathbf{x}_j^i, y_j^i, d_j^i)$ represents the input data of the j -th utterance for domain i , the corresponding speaker label, and the domain label. Domain i has M^i utterances in total. The parameters θ_d minimize the domain classification loss while the parameters θ_s minimize the speaker classification loss. The parameters θ_g minimize the speaker classification loss and simultaneously maximize the domain classification

loss, where the former makes the embeddings discriminative for speaker identity and the latter encourages domain-invariant embeddings to emerge in the course of optimization. The parameter λ is a trade-off factor between the two losses. This process is implemented by a gradient reversal layer (GRL) [26]. The gradient reversal layer has no trainable parameters, and acts as an identity transform in the forward propagation. During the back-propagation, the GRL takes the gradient from the subsequent level, multiplies it by $-\lambda$ and passes it to the preceding layer.

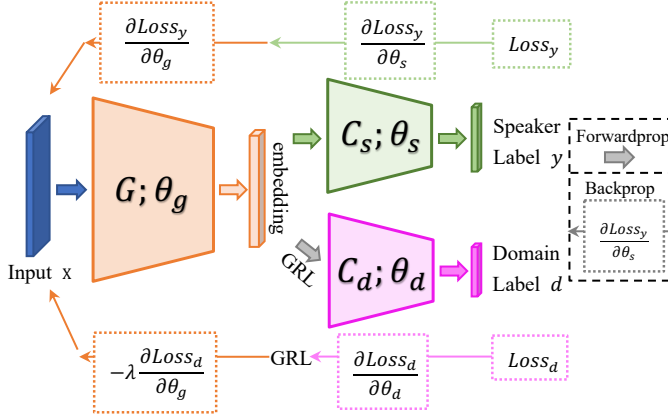


Fig. 1: Domain adversarial training framework.

IV. DISCREPANCY MINIMIZATION ON FEATURES

To mitigate the affect of domain discrepancy, multi-source domain adaptation based on discrepancy minimization, bridging the classification and discrepancy minimization, is employed to extract domain-invariant representations for all domains by means of respectively aligning the distributions of all domain pairs at both the frame-level and segment-level.

A. Pair-wise Distribution Alignment

Maximum mean discrepancy (MMD) is a pair-wise distribution discrepancy measure which is employed over a probability space by computing the mean squared difference of the statistics of samples [50], [51]. Given that two generated distributions are identical, MMD assumes all corresponding statistics are the same and the distribution discrepancy will asymptotically equal 0. The definition in Eq. (3) estimates the maximum mean discrepancy between each domain pair,

$$\mathcal{D}(\mathcal{X}_a, \mathcal{X}_b) = \|\mathbb{E}_{\mathbf{x} \sim \mathcal{X}_a} \Phi(\mathbf{x}^a) - \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_b} \Phi(\mathbf{x}^b)\|^2, \quad (3)$$

where MMD computes the mean square distance between the two collections $\mathbf{X}_a = \{\mathbf{x}_i^a\}_{i=1}^{|\mathcal{X}_a|}$ and $\mathbf{X}_b = \{\mathbf{x}_j^b\}_{j=1}^{|\mathcal{X}_b|}$ of *i.i.d.* sampling from domain \mathcal{X}_a and \mathcal{X}_b , $\Phi(\cdot)$ denotes a feature map, where the mapping is from acoustic frame-level features to an embedding space. Given $|\mathbf{X}_a| = L$ and $|\mathbf{X}_b| = M$, Eq. (3) can be expanded as,

$$\begin{aligned} \mathcal{D}(\mathcal{X}_a, \mathcal{X}_b) &= \frac{1}{L^2} \times \sum_{i=1}^L \sum_{i'=1}^L \Phi(\mathbf{x}_i^a)^\top \Phi(\mathbf{x}_{i'}^a) \\ &\quad - \frac{2}{LM} \times \sum_{i=1}^L \sum_{j=1}^M \Phi(\mathbf{x}_i^a)^\top \Phi(\mathbf{x}_j^b) \\ &\quad + \frac{1}{M^2} \times \sum_{j=1}^M \sum_{j'=1}^M \Phi(\mathbf{x}_j^b)^\top \Phi(\mathbf{x}_{j'}^b). \end{aligned} \quad (4)$$

The dot product can be replaced with the kernel function $k(\cdot, \cdot)$,

$$\begin{aligned} \mathcal{D}(\mathcal{X}_a, \mathcal{X}_b) &= \frac{1}{L^2} \times \sum_{i=1}^L \sum_{i'=1}^L k(\mathbf{x}_i^a, \mathbf{x}_{i'}^a) \\ &\quad - \frac{2}{LM} \times \sum_{i=1}^L \sum_{j=1}^M k(\mathbf{x}_i^a, \mathbf{x}_j^b) \\ &\quad + \frac{1}{M^2} \times \sum_{j=1}^M \sum_{j'=1}^M k(\mathbf{x}_j^b, \mathbf{x}_{j'}^b). \end{aligned} \quad (5)$$

A widely-used kernel function is the radial basis function (RBF) kernel, which ensures that the MMD measure contains all moments of data in the feature space [52]. This kernel function is written as,

$$k(\mathbf{x}^a, \mathbf{x}^b) = \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}^a - \mathbf{x}^b\|^2\right), \quad (6)$$

where σ is a bandwidth parameter of the Gaussian kernel [52].

The motivation for a pair-wise distribution alignment is to ensure that the network predictions are consistent even if inputs are subject to an intrinsic/extrinsic [1] domain shift. As noted in [23], larger domain discrepancy gaps typically exist in deeper layers such as the fully-connected layer generating the embeddings. Therefore, we will minimize the discrepancy among embeddings produced by data from each domain. Additionally, the employed network uses an adaptive training scheme where samples are grouped into short segments (400 frames randomly out of 1000 frames) in each mini-batch, similar to the pre-training scheme based on sampling arbitrary-length segments from 200 to 400 frames. Moreover, the statistics pooling may also distort speech sequence features, especially based on temporal-related information. To alleviate this possible embedding distribution shift caused by speech sampling and statistics pooling, we also adapt frame-level features including the last TDNN layer's output before statistics pooling to avoid this inaccurate discrepancy estimate. Let $\mathcal{O}_a^l = \{\mathcal{O}_i^l\}_{i=1}^{|\mathcal{X}_a^l|}$ denote the collection of l -layer outputs from the distribution \mathcal{O}_a^l for domain \mathcal{X}_a . Multiple domains will possess a domain shift with each other, where domain-invariant representations for each paired domain can be learned by minimizing the $Loss_{mmd}$ as follows,

$$\begin{aligned} Loss_{mmd} &= \binom{N}{2}^{-1} \sum_{i=1}^{N-1} \sum_{j=i+1}^N \left(\mathcal{D}(\mathcal{O}_i^{\mathcal{F}_5}, \mathcal{O}_j^{\mathcal{F}_5}) \right. \\ &\quad \left. + \mathcal{D}(\mathcal{O}_i^{f_{c1}}, \mathcal{O}_j^{f_{c1}}) \right), \end{aligned} \quad (7)$$

where \mathcal{F}_5 is the last TDNN layer before statistics pooling in the universal embedding extractor, and f_{c1} is the first fully-connected layer which generates the embeddings. The discrepancy loss is computed between each pair of N domains.

B. Domain-specific Classification

Each domain-specific subnet is followed by a softmax classifier. We use a softmax cross-entropy \mathcal{J} loss for each classifier to ensure that the embedding distribution performance is improved for each domain. The classification loss function is written as,

$$Loss_{cls} = \sum_{i=1}^N \sum_{j=1}^{M^i} \mathcal{J}(C_i(G(\mathbf{x}_j^i)), y_j^i). \quad (8)$$

Given N domains, classification loss is computed for each domain-specific subnet. Here, (\mathbf{x}^i, y^i) represents the acoustic frame-level input feature of an utterance for domain i and the corresponding speaker label. Domain i has M^i utterances in total, and G represents the universal embedding extractor, which maps the input feature to a universal embedding. Finally, C_i is the classification subnet of the domain i out of N domains after employing the embedding extractor.

C. Discrepancy Minimization Adaptation Framework

Given an input data sequence, the universal feature extractor projects the sequence data into a temporal orderless embedding. The classification component will have four independent subnets corresponding to specific domains. The framework of the cross-domain adaptation in each step is illustrated in Fig. 2. Here, the multi-task loss function is formulated as,

$$Loss_{total} = \mu(Loss_{mmd}) + Loss_{cls}, \quad (9)$$

where μ is a variant adaptation factor with a progressive schedule from 0 to 1 in order to stabilize parameter sensitivity in the early adaptation stage.

Here, we employ our pre-trained x-Vector system as the universal embedding extractor, which is extended with classification subnets. The distributions of each domain are aligned simultaneously by minimizing domain discrepancies. Subsequently, the domain-invariant representations are specifically learned. Furthermore, the domain-specific classification subnets are employed to optimize recognition performance for individual domains of interest.

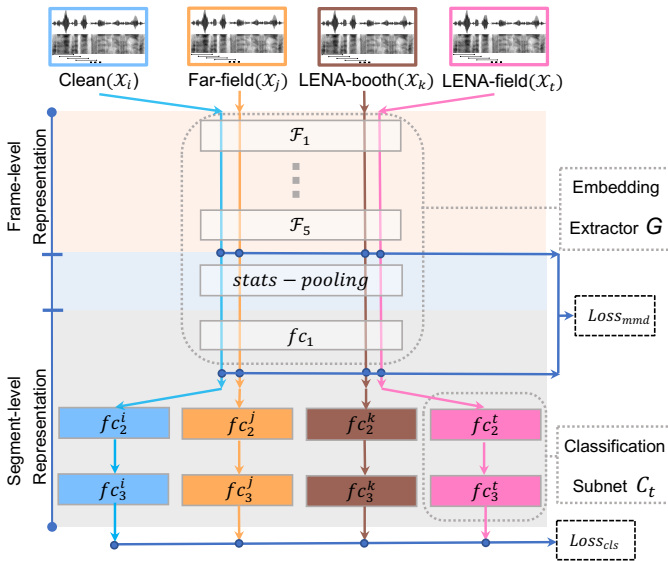


Fig. 2: Discrepancy minimization adaptation framework.

V. MOMENT MATCHING ON THE CLASSIFICATION POSTERIOR

For this study, four domains of our CRSS-Forensics dataset are investigated (details later in Sec. V). In general, multi-source domain data reduces the effectiveness of any single domain adaptation method. Additionally, the domain discrepancy

will also vary in each pair of the domain-specific data. For this corpus, Clean, Far-field, and LENA-booth data were collected in the same recording environment (sound booth). For LENA-field data, speech data was collected using a portable LENA recording unit worn by the participant, with recording environments including seven pre-defined indoor and outdoor locations. Based on these corpus specifics, data collected in the sound booth have marginal domain shifts among each dataset (e.g. only close-talk mic (CTM) vs. desktop or distance mics at 4 ft and 8 ft). In contrast, a distinct domain discrepancy exists between LENA-field data and all data collected in the environment controlled sound booth. Therefore, we consider a multi-source domain adaptation based on moment matching and disposed multiple complex adversarial training procedures [41]. This method is employed to minimize the inter-domain discrepancies and transfer knowledge learned from sound-booth data to the more diverse LENA-field data by dynamically aligning moments of their feature distributions.

A. Moment-matching Components

Given labeled data collections $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$ with distinct domain discrepancies for input feature space \mathcal{X} , the moment-matching process aims to find discriminative representations in the hypothesis embedding space \mathcal{H} , which minimizes the testing error on each domain. In [41], the domain discrepancy was measured with the Moment Distance. Here, we use the measurement in Eq. (3) instead to define the moment distance. The kernel function lifts the sample vectors into an infinite dimensional feature space and covers all orders of statistics, consequently minimizing MMD with this kernel which is equivalent to minimizing a distance between all moments of the two distributions [50].

We employ a moment-matching model, which comprises of a universal embedding extractor G , along with a set of N classifiers $\mathcal{C} = \{C_1, C_2, \dots, C_N\}$ (e.g. this is the same setup as described in Sec. III). The MMD measurement minimizes the moment-related distance between domains as defined in Eq. (3). The overall loss function in Eq. (9) is therefore rewritten as the following objective function,

$$\min_{G, \mathcal{C}} \sum_{i=1}^N \mathcal{J}_{\mathcal{X}_i} + \mu \min_G \binom{N}{2}^{-1} \sum_{i=1}^{N-1} \sum_{j=i+1}^N (\mathcal{D}(\mathcal{O}_i^{\mathcal{F}_5}, \mathcal{O}_j^{\mathcal{F}_5}) + \mathcal{D}(\mathcal{O}_i^{\mathcal{F}_{C_1}}, \mathcal{O}_j^{\mathcal{F}_{C_1}})), \quad (10)$$

where $\mathcal{J}_{\mathcal{X}_i}$ is a softmax cross-entropy loss for the classifier C_i for domain \mathcal{X}_i , and μ is a trade-off parameter with a progressive schedule. The objective of moment-matching adaptation is to match different distributions by minimizing the moment distance between multiple acoustic domains. Furthermore, for our task, we also intend to leverage knowledge learned from the noise-free sound-booth data to recalibrate the distribution for our diverse LENA-field data.

B. Adversarial Training Schema

We follow the training paradigm suggested in [27], in order to utilize the domain-specific decision boundaries. Considering the relationship between class boundary and LENA-field samples, the paired domain-specific classifiers are taken as

a discriminator to detect the presence of LENA-field samples from that reflecting the sound-booth domain. Paired classifiers are likely to classify those outliers in LENA-field samples differently. The tandem adaptive training includes the following three steps:

1) Train the universal embedding extractor G and classifier collection \mathcal{C} to minimize moment distances as in Eq. (3) among domains and perform classification on each domain. Model parameters are updated using the objective from Eq. (10).

2) Fix the parameters of G , so as to maximize the discrepancies of classifier pairs. To measure the discrepancy of the two classifiers, we utilize the MMD as in Eq. (3), which formulates the objective in this training step,

$$\sum_{i=1}^{N-1} (\min_{C_i} \mathcal{J}_{\mathcal{X}_i} - \mathcal{D}(C_i(\mathbf{X}_i), C_N(\mathbf{X}_N))) + \min_{C_N} (\mathcal{J}_{\mathcal{X}_N} - \frac{1}{N-1} \sum_{j=1}^{N-1} \mathcal{D}(C_j(\mathbf{X}_j), C_N(\mathbf{X}_N))), \quad (11)$$

where $C_i(\mathbf{X}_i)$ and $C_N(\mathbf{X}_N)$ represent the probability outputs of C_i and C_N respectively from one of the sound-booth domains and LENA-field domains. The classification loss on each domain is added to stabilize system performance.

3) Finally, we fix \mathcal{C} and train G to minimize the discrepancy of each classifier pair. The objective of this step is written as,

$$\min_G \sum_{i=1}^{N-1} \mathcal{D}(C_i(\mathbf{X}_i), C_N(\mathbf{X}_N)). \quad (12)$$

This entire procedure is summarized as Algorithm 1. For this solution, we train the classifiers and generator in an adversarial manner until the entire network (see in Fig. 3) reach a point of convergence.

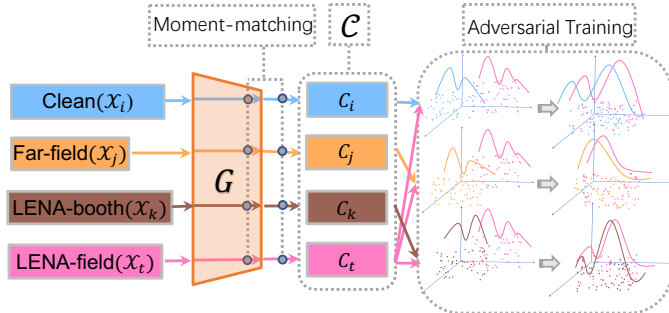


Fig. 3: Moment-matching adaptation architecture.

VI. EXPERIMENT

A. Data Description

1) **VoxCeleb**: We use the Vox2 and Vox1 dev corpora for embedding training [48], which is extracted from videos based on YouTube as training data for our pre-trained systems. Videos included in the dataset are recorded in a large number of challenging visual and auditory environments, including background conversations, laughter, overlapping speech, and varying room acoustics. Over 2.2 million utterances from ≈ 7300 speaker identities were used with corresponding annotation for speaker labels. Following a baseline Kaldi recipe, we use the dev and test splits from Vox2 and the dev split from Vox1 for embedding-oriented pre-training.

Algorithm 1 Moment-Matching Adaptation Network

Input: $\mathcal{X} = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N\}$, pre-trained G
Output: pre-trained G and a set of N classifiers $\mathcal{C} = \{C_1, C_2, \dots, C_N\}$

- 1: Given T_1 and T_2 training iterations
- 2: **for** t in $1 : T_1$ **do**
- 3: **for** j in $1 : N$ **do**
- 4: Sample mini-batch $\{\mathbf{x}_i^j, \mathbf{y}_i^j\}_{i=1}^m$ from \mathbf{X}_j
- 5: Feed $\{\mathbf{x}_i^j, \mathbf{y}_i^j\}_{i=1}^m$ to G
- 6: Feed embeddings from G to C_j
- 7: **end for**
- 8: **Update** G and \mathcal{C} according to Equation(10)
- 9: **end for**
- 10: **for** t in $1 : T_2$ **do**
- 11: **for** j in $1 : N$ **do**
- 12: Sample mini-batch $\{\mathbf{x}_i^j, \mathbf{y}_i^j\}_{i=1}^m$ from \mathbf{X}_j
- 13: Feed $\{\mathbf{x}_i^j, \mathbf{y}_i^j\}_{i=1}^m$ to G
- 14: Feed embeddings from G to C_j
- 15: **end for**
- 16: Fix G
- 17: **Update** \mathcal{C} according to Equation(11)
- 18: Fix \mathcal{C}
- 19: **for** t' in $1 : 4$ **do**
- 20: **Update** G according to Equation(12)
- 21: **end for**
- 22: **end for**

2) **CRSS-Forensic**: As noted earlier, the CRSS-Forensics corpus¹ contains read speech, prompted speech, and spontaneous speech in three conditions: clean (noise-free recorded in the sound booth), field recordings (with background noise and reverberation), and high stress (actual police interviews). Two phases are included in the recording process. Phase-1 contains speech data recorded in a controlled noise-free sound booth, and diverse public acoustic field environments. Speech in Phase-2 is collected in a law enforcement facility, using an interview room with an actual police officer/detective. Figure 5 shows sample recording environments for the noise-free sound booth, public field environments, and police interview room.

For the sound booth in Phase-1, speech data is simultaneously recorded using multiple wired/wireless microphones (sample rate: 44.1 kHz) and a participant body-worn mobile data collection platform called LENA unit (sample rate: 16 kHz). Microphones are positioned at 4 different locations in a $13' \times 13'$ sound booth (SB); one being a close talk microphone (CTM), the other three representing distances from each microphone to the speaker are 18-24 in (Desktop Mic (DTM)), 4 ft, and 8 ft, as shown in Fig. 4. In LENA-field environments, speech is collected by a LENA unit worn by the participant, with recording environments that include seven indoor and outdoor locations (7IOL): office, hallway, cafeteria, parking lot, game room, lobby, walking path (see in Fig.5 (a-g)). In Phase-2, speech data is simultaneously recorded using a participant body-worn LENA unit and a microphone. Here a detective conducts an investigative interview of the participant concerning a specific scenario while following standard procedures in law enforcement interview room (LEIR) (see in Fig. 5 (h)).

Table I summarizes the specific acoustic data size for each session. For the 75 speakers in the corpus, 65 are native English speakers and 10 non-native speakers, with 27 male

¹The CRSS-Forensics corpus will be released with a license.

TABLE I
DATA STATISTICS FOR CRSS-FORENSICS CORPUS

Info	Session Name	Duration	Speaker
Phase-1.SB	Clean (CTM & DTM)	32 h/channel	75
	Far-field (4 ft & 8 ft)	32 h/channel	75
	LENA-booth	33.9 h	75
Phase-1.7IOL	LENA-field	99.4 h	75
Phase-2.LEIR	Interview (LENA & Mic)	20.4 h	58

speakers and 48 female speakers. Each participant was allowed to opt-out of Phase-2 (i.e., IRB protocol due to high-stress level exposure), so there are 17 speakers absent from the Phase-2 police interview set.

In this study, data from Phase-1 is used for multi-source domain adaptation. We note that various recording environments are considered as the extrinsic characteristics for audio samples, while speech from speakers under stress for the Interview in Phase-2 contain intrinsic variations. Consequently, data in Phase-2 is not compatible with environmental mismatch data in Phase-1 for domain-invariant information extraction. We consider 16 speakers (8 male, 8 female) out of 75 speakers for each set to perform evaluation. There exists no speaker overlap between the training and evaluation sets, abiding by an open-set protocol. The number of test trials is over 22,000 total. For sound booth data, speech data with CTM and DTM are designated as the Clean set; data collected from the remaining two distant mics (4 ft & 8 ft) are taken as Far-field data; and data recorded by the LENA body-worn unit was used to explore channel mismatch influence.

B. Pre-training System Setup

Gender independent i-Vector extractors were trained on the VoxCeleb dataset to produce 400-dimensional i-Vectors. 20-dimensional MFCCs were then augmented with their delta and double-delta coefficients, producing a set of 60-dimensional MFCC feature vectors.

In order to implement a competitive and fair baseline, we developed the x-Vector system. Our model is similar to an x-Vector Kaldi recipe² with respect to VoxCeleb corpora and network architecture. The model architecture consists of 5 time-delay layers, which model temporal context information, followed by a statistical pooling layer to map into a fixed-dimensional vector at the segment-level. This is followed by two fully-connected layers with 512 units in each layer and the probability output layer. We extract 30-dimensional MFCC features using a frame width of 25ms and window shift of 10ms. Training data is augmented with noise, music, and babble speech from the MUSAN corpus [53], and reverberation of the RIR NOISES³ corpus. The augmented data consist of 7323 speakers and 2.2M utterances. Specially, all utterances shorter than 4 seconds in duration, and all speakers with fewer than 8 utterances are set aside in the data pre-processing phase. Cepstral mean normalization with a sliding window of 3 seconds was employed to suppress channel effects. We use an Adam optimizer with betas of (0.9, 0.98) to update model

parameters, initializing the learning rate of 1e-3. The learning rate was adjusted with the warm-up scheduling named “Noam” in [54]. Batch normalization and Dropout are also used to perform regularization at each layer. Finally, a mini-batch of 32 samples is used at each iteration.

C. Fine-tuning Setup

In this work, all adaptation methods only update parameters of the last time-delay layer before statistics pooling and the first fully-connected layer after pooling in the pre-trained x-Vector model. The last fully-connected layer is replaced according to speaker labels from the CRSS-Forensic data. We perform fine-tuning on the pre-trained model with our data using the Adam optimizer to retrain the model for 40 epochs using a batchsize of 64. The learning rate is scheduled using the formula,

$$\eta_p = \frac{\eta_0}{(1 + \alpha p)^\beta}, \quad (13)$$

where $\eta_0 = 1e-4$, $\alpha = 10$, $\beta = 0.75$ and p is set to linearly increase from 0 to 1 corresponding to the training steps.

D. Domain Adversarial Training System Setup

We keep the pre-trained x-Vector model as the feature extractor, and extract embeddings from the first fully-connected layer after statistic pooling. For the speaker label classifier, we retain the three fully-connected layers ($Embedding \rightarrow 512 \rightarrow 512 \rightarrow 59$), and use a simpler architecture ($Embedding \rightarrow GRL$ (described in Sec. III) $\rightarrow 128 \rightarrow 4$) for domain classification. The model is trained on 64-sized batches. In order to suppress noisy signals from the domain classifier at the early training stages instead of fixing the trade-off factor λ , we gradually change this value from 0 to 1 using the following schedule:

$$\mu = \frac{2}{1 + \exp(-\theta p)} - 1, \quad (14)$$

where $\theta = 10$, p is set to linearly increase from 0 to 1 corresponding to the training steps.

E. Discrepancy-minimizing System Setup

In addition to the pre-trained model, we keep the first fully-connected layer after statistics pooling as part of the embedding extractor. Embeddings from each domain are processed by a fully-connected layer with 512 units and the final layer which outputs speaker posterior probabilities, respectively. The discrepancy-minimizing model is trained on the CRSS-Forensic data for 40 epochs using a batchsize of 64. We use the Adam optimizer to update parameters of the partial pre-trained model with a learning rate of 1e-4, and for each classification subnet, the learning rate is set to 1e-3. Since there exists no parameter-wise differences between each subnet in the early adaptation stage, Eq.(9) may result in noisy activations. To stabilize parameter sensitivity, a progress strategy [26] is used for Eq.(9) as noted in Eq.(14).

²<https://kaldi-asr.org/models/m7>

³<http://www.openslr.org/28>

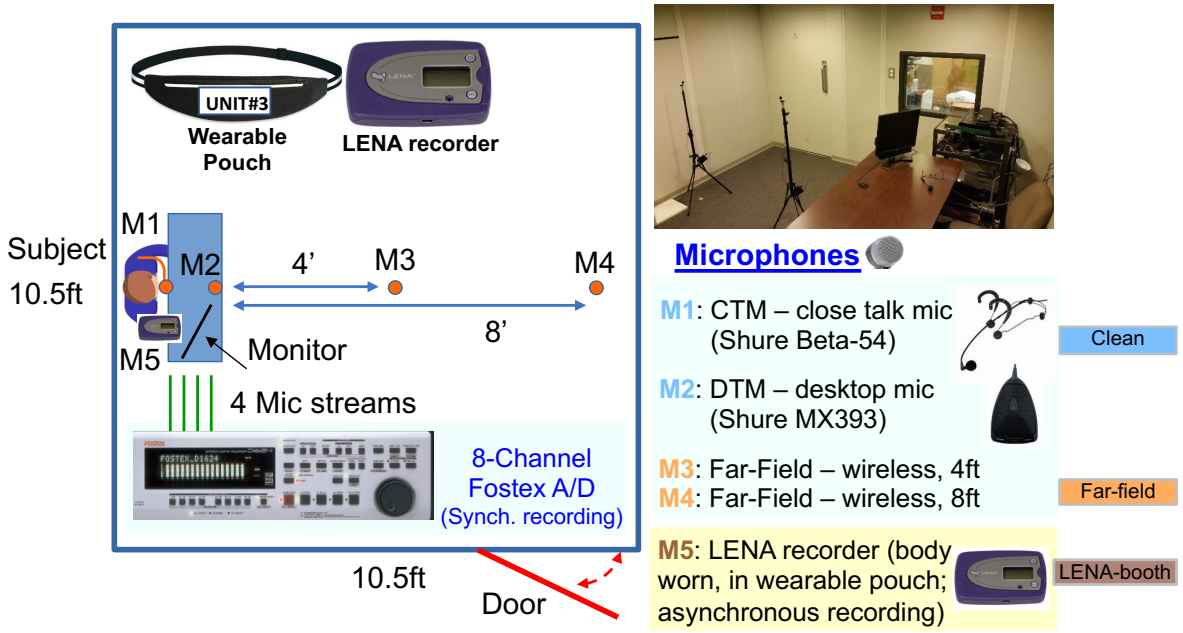


Fig. 4: Sound booth forensic voice data collection setup, including 5 audio streams (M1: close-talk mic, M2: desk-top mic, M3 & M4: far-field distance mics, M5: asynchronous body-worn LENA recorder).

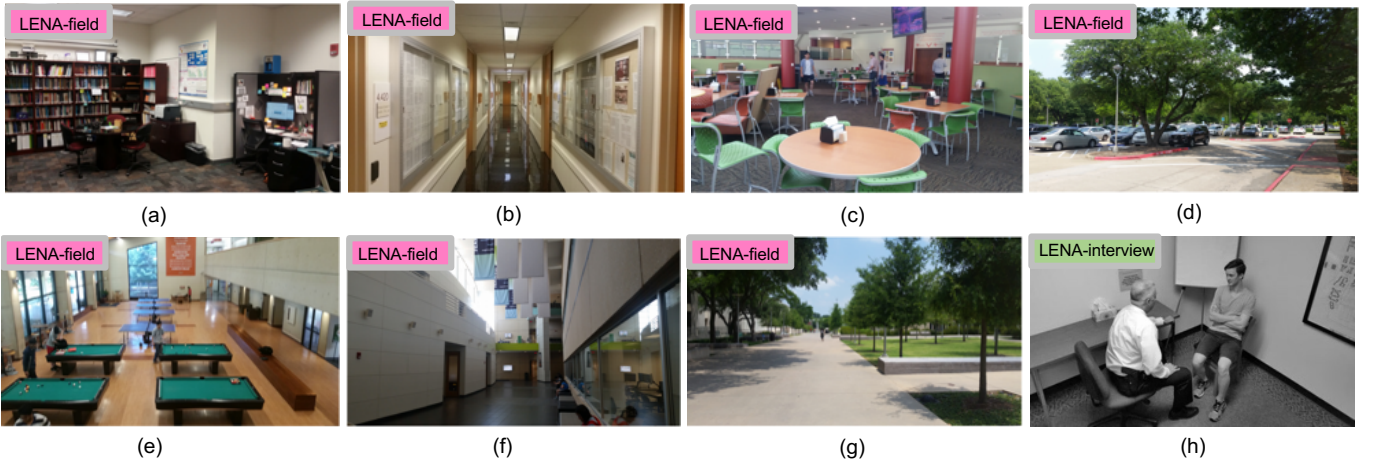


Fig. 5: LENA-Field forensic audio collection locations (a) to (g); plus Police Interview location (LENA-Police) recording environments (h).

F. Moment-matching System Setup

The model setup for moment-matching system is basically consistent with the previous description for the discrepancy-minimizing system with the exception for the training schema and loss functions constituted by another discrepancy measurement (details in Sec. V). The moment-matching model is trained for 40 epochs as T_1 in Algorithm 1 then proceeds for another 40 epochs as T_2 in Algorithm 1 with a batchsize of 64.

For each system, we take embeddings from the outputs of the first fully-connected layer after statistics pooling for evaluation purposes, and score trials using PLDA [55] after performing dimensionality reduction to 200 using LDA and length-normalization. Here, LDA and PLDA multi-conditional training are conducted in each system with generated embeddings from the training portion in CRSS-Forensic corpus containing 59 speakers, which can compensate for domain mismatch.

VII. RESULT AND ANALYSIS

This section focuses on the analysis of each system implemented based on setups described in Sec. VI. To evaluate these experiments, we adopt several measurement criteria concentrating on evaluating discrimination abilities and calibration of speaker recognition systems. In terms of speaker recognition system evaluation, the trade-off between missed speakers P_{miss} and false alarms P_{FA} has always been a key diagnostic tool. The Detection Error Trade-off (DET) curve [56] reflects what happens as the decision threshold is swept across the entire operating range. Noting that P_{miss} and P_{FA} move in opposite directions as the decision threshold is shifted, a point where $P_{miss} = P_{FA}$ called the Equal Error Rate (EER), provides a standard point for the discrimination capability of the system. However, the EER does not measure calibration (the ability to set decision thresholds). It is noted that the recognition system actually produces the log-likelihood-ratio

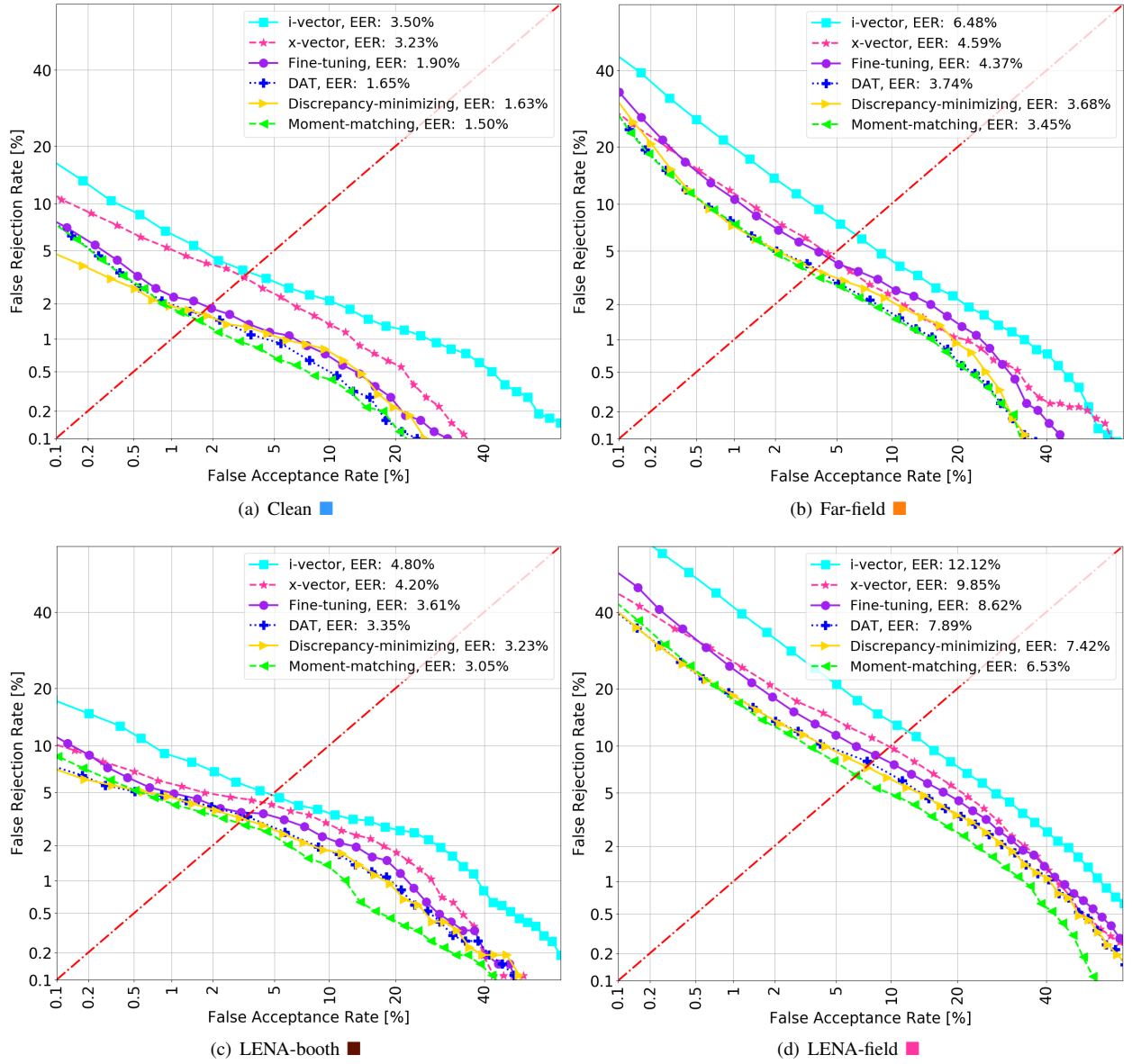


Fig. 6: DET curves for each system on four domains.

\mathcal{L}_t of the score for each trial. The Log-likelihood-ratio cost function [57] C_{llr} is a simultaneous measurement of the discrimination abilities of the log-likelihood-ratio scores and calibration for application-independent detection decisions, which is formulated as,

$$C_{llr}(\mathcal{L}_t) = \frac{1}{2} \left(\frac{1}{N_{tar}} \sum_{t \in T_{tar}} \log(1 + e^{-\mathcal{L}_t}) + \frac{1}{N_{non}} \sum_{t \in T_{non}} \log(1 + e^{\mathcal{L}_t}) \right), \quad (15)$$

where \mathcal{L}_t is the log-likelihood-ratio score of trial t , T_{tar} is a set of N_{tar} target trials and T_{non} is a set of N_{non} non-target trials. Furthermore, C_{llr} can be minimized as measured on a warped log-likelihood-ratio score $\mathcal{L}'_t = w(\mathcal{L}_t)$ scale using a monotonic rising warping function w , resulting in the performance measure $C_{llr}^{min} = C_{llr}(\{w(\mathcal{L}_t)\})$. Finally, w can be optimized using the Pool Adjacent Violators (PAV) approach [57]. We visualized each speaker recognition system performance on each acoustic domain using DET curves (see

in Fig. 6).

A. Score Calibration

The recognition system produces the log likelihood ratio in terms of PLDA scores, but the scores are uncalibrated which may adversely affect the validity and reliability of this evaluation. Score calibration has been recognized as an important component in effective evaluation of current speaker recognition systems [57]–[59]. Thus, we calibrate the log-likelihood-ratio (LLR) scores by finding a linear transform that optimizes the CLLR measure to reach a value of C_{llr}' closer to C_{llr}^{min} . In this study, we employ a commonly-used linear calibration transformation,

$$s' = w_0 + w_1 s, \quad (16)$$

where an uncalibrated score s is transformed into a calibrated score s' using offset w_0 and scaling factor w_1 parameters. Logistic regression optimization [60] is employed to acquire

TABLE II
SCORE CALIBRATION RESULT FOR EACH SYSTEM ON FOUR DOMAINS

	Clean			Far-field			LENA-booth			LENA-field		
	C_{llr}	C'_{llr}	C_{llr}^{min}	C_{llr}	C'_{llr}	C_{llr}^{min}	C_{llr}	C'_{llr}	C_{llr}^{min}	C_{llr}	C'_{llr}	C_{llr}^{min}
i-Vector	0.241	0.172	0.138	0.335	0.247	0.231	0.364	0.234	0.176	0.591	0.428	0.399
x-Vector	0.211	0.130	0.110	0.323	0.189	0.168	0.331	0.180	0.136	0.610	0.315	0.314
Fine-tuning	0.200	0.098	0.072	0.364	0.183	0.167	0.288	0.144	0.123	0.621	0.323	0.293
DAT	0.272	0.118	0.064	0.312	0.145	0.134	0.294	0.114	0.105	0.528	0.276	0.254
Discrepancy-minimizing	0.172	0.085	0.063	0.298	0.155	0.139	0.231	0.124	0.109	0.526	0.284	0.256
Moment-matching	0.272	0.080	0.060	0.276	0.140	0.131	0.211	0.111	0.100	0.479	0.254	0.233

the two calibration parameters w_0 and w_1 . We summarized score calibration results for each speaker recognition system on each acoustic domain using DET curves (as shown in Table II), where C'_{llr} corresponds to C_{llr} after calibration.

B. Location Analysis

Next, speaker recognition (SR) performance is assessed in terms of EER and C_{llr} for both i-Vector and x-Vector system over the range of evaluation datasets (as shown in Table III). The pre-trained x-Vector system shows better speaker recognition performance, so it is taken as the embedding extractor for subsequent adaptation methods. Additionally, it is noted that the x-Vector architecture is based on a deep neural network, which allows for fine-tuning and to concatenate with other deep-learning structures. In terms of impact due to domain mismatch on system performance, channel, speaker-to-mic distance, and environmental noise all exert some mismatch influence on system recognition performance, with noise mismatch having the greatest impact.

TABLE III
SR RESULT FOR I-VECTOR & X-VECTOR SYSTEMS

	i-Vector			x-Vector		
	EER	C'_{llr}	C_{llr}^{min}	EER	C'_{llr}	C_{llr}^{min}
Clean	3.50%	0.172	0.138	3.23%	0.130	0.110
Far-field	6.48%	0.247	0.231	4.59%	0.189	0.168
LENA-booth	4.80%	0.234	0.176	4.20%	0.180	0.136
LENA-field	12.12%	0.428	0.399	9.85%	0.315	0.314

TABLE IV
SR RESULT FOR X-VECTOR SYSTEM IN 7 LOCATIONS

7IOL	EER	C_{llr}	C'_{llr}	C_{llr}^{min}
Cafeteria	12.61%	0.807	0.409	0.372
Game room	13.78%	0.900	0.444	0.416
Hallway	9.23%	0.632	0.314	0.292
Lobby	8.55%	0.493	0.293	0.267
Office	6.64%	0.395	0.241	0.224
Parking Lot	7.35%	0.401	0.262	0.240
Walking Path	9.01%	0.545	0.311	0.290

The LENA-field set includes 7 naturalistic locations (7IOL as shown in Fig. 5). We evaluate x-Vector system performance on data across each environmental location (as shown in Table IV). Results show that speech data captured in public cafeteria and game room locations had the lowest speaker recognition results versus other locations. Cafeteria and game room data

contain secondary people talking, and sporadic background music and random noise/sound events which are clearly heard especially in the game room. Speaker identity is more easily discriminated with data from the office context, since noise content is less, though background talking can occur at times. Other locations contain varying amounts of reverberation and ambient noise, also resulting in degradation in recognition performance.

C. Fine-tuning Layers Selection

In order to explore the best fine-tuning result, we performed fine-tuning of the pre-trained x-Vector system for different layers. Here, we present the fine-tuning results across 3 different options: ($\mathcal{F}_4, \mathcal{F}_5, f_{c1}$), (\mathcal{F}_5, f_{c1}) and f_{c1} (definition see Fig. 2) as shown in Table V.

TABLE V
 C_{llr} RESULT FOR X-VECTOR MODEL FINE-TUNING

	$\mathcal{F}_4, \mathcal{F}_5, f_{c1}$		\mathcal{F}_5, f_{c1}		f_{c1}	
	C'_{llr}	C_{llr}^{min}	C'_{llr}	C_{llr}^{min}	C'_{llr}	C_{llr}^{min}
Clean	0.102	0.075	0.098	0.072	0.100	0.074
Far-field	0.174	0.191	0.183	0.167	0.187	0.170
LENA-booth	0.150	0.127	0.144	0.123	0.147	0.125
LENA-field	0.334	0.304	0.323	0.293	0.328	0.298

Table V shows that fine-tuning of the pre-trained x-Vector model achieves the best result for C'_{llr} and C_{llr}^{min} when applied in the last layer before statistics pooling, and the first layer after pooling. By performing fine-tuning in the proper layers, knowledge of the pre-training data is effectively transferred towards the current model, and the model also learns effective speaker information for the new dataset. We only fine-tune subsequent layers to maintain the learned universal speaker features from undue distortion. The fine-tuned system lowers pre-trained system's EER with a relative decrease of 41.18%, 4.79%, 12.49%, 14.05% in each set, respectively, with an averaged relative decrease in EER of 18.13%. Obviously, the Clean set benefits the most from fine-tuning.

D. Adaptive Training with Domain Information

Fine-tuning effectively improves speaker recognition performance for the pre-trained x-Vector system, though the achieved improvement is unbalanced across each domain. Therefore, it is necessary to explore other options to improve use of domain information for better speaker recognition performance. This can be achieved by DAT [49] which uses a

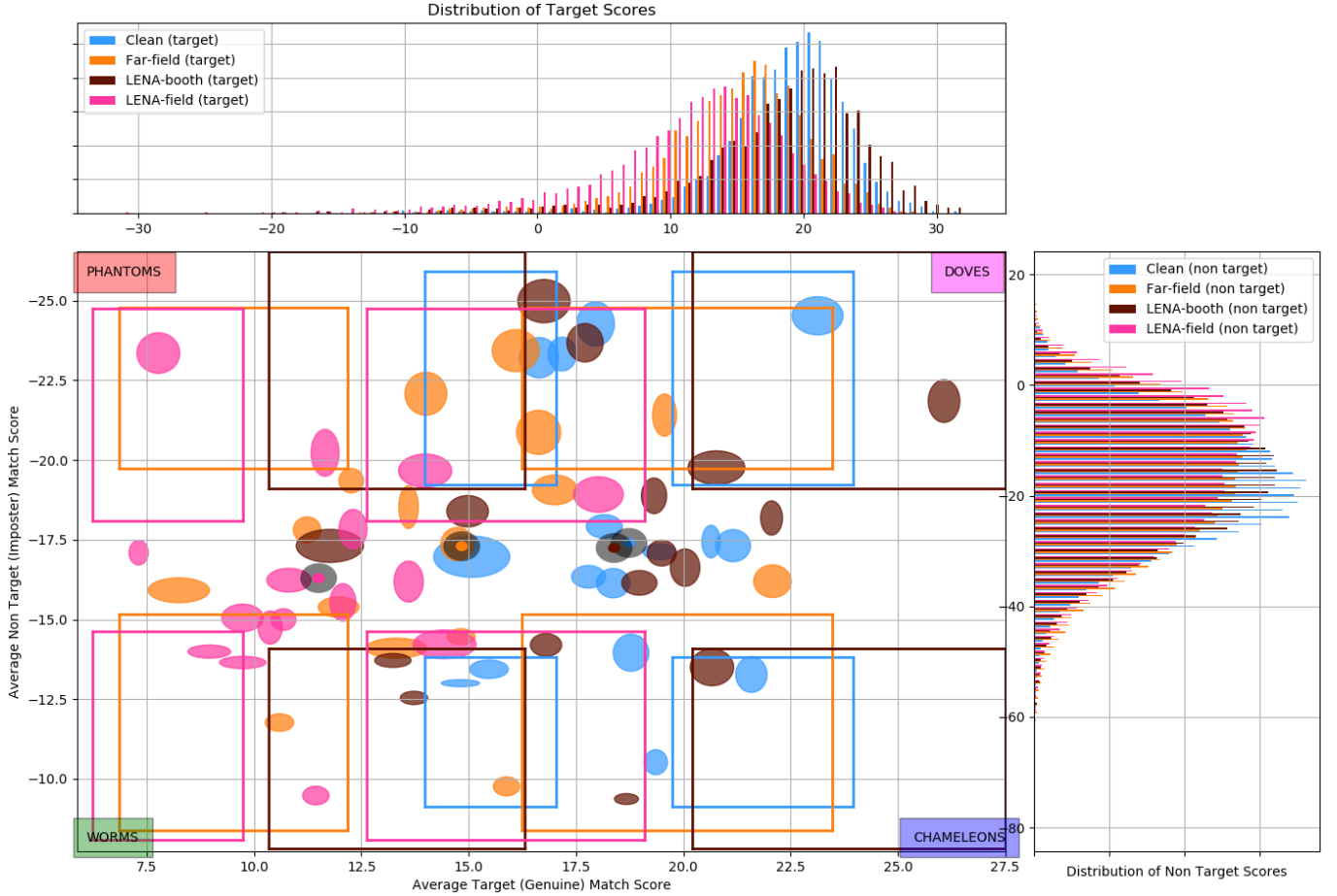


Fig. 7: Zoo-plot for discrepancy minimization system.

gradient reversal layer to remove the domain variation and projects the alternate domain data towards the same subspace. This approach learns a domain-invariant and speaker-discriminative feature representation. The DAT system lowers EER of the pre-trained system with a relative decrease of 48.92%, 18.52%, 20.24%, 19.90% for each set, respectively, with an averaged relative EER decrease of 26.90%.

E. Frame-level Effectiveness in Discrepancy Minimization

An alternative strategy to address mismatch is to consider a multi-task loss function in order to improve the pre-trained system performance across acoustic environments. The main objective is to minimize the domain discrepancy so as to achieve system performance gains for each domain.

TABLE VI
SR RESULT FOR DISCREPANCY MINIMIZATION SYSTEM

	frame-level			frame-level & seg-level		
	EER	C'_{llr}	C_{llr}^{min}	EER	C'_{llr}	C_{llr}^{min}
Clean	1.65%	0.081	0.060	1.63%	0.083	0.063
Far-field	3.84%	0.156	0.141	3.68%	0.151	0.139
LENA-booth	3.27%	0.123	0.107	3.23%	0.127	0.109
LENA-field	7.50%	0.280	0.254	7.42%	0.289	0.256

Based on this motivation, we consider discrepancy minimization within an adaptation procedure at both the segment-level and frame-level to avoid inaccurate discrepancy esti-

mation caused by the domain-wise embedding distribution deviation. Results from Table VI show improvement with adaptation based on discrepancy minimization, where frame-level adaptation contributes to improvement with a slight EER reduction. Additionally, we utilize a zoo-plot visualization [61] to explore a sample analysis on individual speakers, or speaker groups (see in Fig. 7). The zoo-plot shows a scatter type visualization based on mean values of both target and non-target scores for each speaker label, and speakers who fall within the four quadrants are assigned to animal groups (worms, chameleons, doves, and phantoms) with each set showing different characteristics. The black ellipses show mean values of all target and non-target scores for each domain, with the domain-index color in the center of each eclipse. Speakers toward the upper right corner have lower genuine variability and higher imposter variability. For example, speakers in the CRSS-Forensic LENA-field (recorded in 7 indoor and outdoor locations) tend to be more difficult to verify correctly than those in other domains. This visualization helps reveal potential algorithmic weaknesses against certain classes of speakers and domains. In terms of a statistics comparison, discrepancy minimization does improve EER of the pre-trained system with a relative decrease of 49.54%, 19.83%, 24.67%, 23.09% in each set, respectively. The average relative EER decrease is 29.28%.

F. Reducing the Noise Mismatch Impact

For the CRSS-Forensic corpus, the LENA-field portion consists of 7 diverse individual locations, which differ from the other three forensic portions obtained in the sound booth. As shown in the experiment results from Table IV, the noise mismatch exerts a significant impact on speaker recognition performance. Therefore, here we employ the moment-matching method with an adversarial training schema for adaptation to minimize both the domain shift and simultaneously mitigate the impact of noise.

TABLE VII
SR RESULT FOR MOMENT-MATCHING SYSTEM

	EER	C_{llr}	C'_{llr}	C_{llr}^{min}
Clean	1.50%	0.173	0.08	0.060
Far-field	3.45%	0.276	0.14	0.131
LENA-booth	3.05%	0.211	0.111	0.100
LENA-field	6.53%	0.479	0.254	0.233

The moment-matching advancement is shown to reduce EER of the pre-trained system with a relative decrease of 53.56%, 24.84%, 33.71%, 27.38% in each set, respectively. The average relative EER decrease is 34.87%. As shown in Table VII, moment-matching with adversarial training improves SV performance for both sound-booth and LENA-field datasets by dynamically aligning the distribution of the LENA-field set with sound-booth captured audio sets. This distribution alignment is a mutual recalibration process, which suggests sound-booth data provides more distinguishable speaker information versus LENA-field data. Here, LENA-field data increases data diversity of sound-booth data for better generalization. Specifically, Table VIII summarizes the statistics for system SV results for each location of the LENA-field set. A comparison of Table VIII versus baseline x-Vector result (as shown in Table IV) confirms the dramatic benefits of the proposed solution.

TABLE VIII
SV RESULT FOR THE MOMENT-MATCHING SYSTEM IN 7 LOCATIONS

7IOL	EER	C_{llr}	C'_{llr}	C_{llr}^{min}
Cafeteria	9.48%	0.693	0.334	0.301
Game room	9.36%	0.702	0.358	0.324
Hallway	6.13%	0.491	0.248	0.221
Lobby	5.17%	0.377	0.212	0.181
Office	3.30%	0.254	0.132	0.118
Parking Lot	4.19%	0.279	0.177	0.152
Walking Path	6.18%	0.451	0.231	0.206

To visualize the effect of the moment-matching system on speaker recognition performance for LENA-field data, we further assess the quality of the learned speaker features using a t-distributed Stochastic Neighbor Embedding (t-SNE) plot [62] (see in Fig. 8). We plot embeddings after LDA from 16 speakers of the CRSS-Forensic test set, which are generated by the discrepancy-minimizing system and moment-matching system. In the center of Fig. 8 (a), there is a cluster of outlier samples for the LENA-field speaker embeddings from the discrepancy-minimizing system, which confirms that they are easily misclassified with ambiguous identities. Alternatively,

Fig. 8 (b) shows a sparse confusion cluster in the center which highlights how utilizing the domain-specific decision boundaries (noted in Sec. VI.B) works to improve speaker discrimination of outlier samples near the classification boundaries by dynamically aligning distributions in the moment-matching system. Several previous outlier samples in Fig. 8 (a) are also reclustered into corresponding groups in Fig. 8 (b), where most of the remaining samples could be actual outliers such as noise and speech of non-target speakers. Speech in the LENA-field will often contain sporadic noise and non-related speech due to the naturalistic field locations, which are also labeled as target speakers with a coarse-grained transcription.

VIII. FORENSIC SPEAKER RECOGNITION

The object of forensic speaker recognition is to assist in the "trier of fact" (i.e., a judge, a panel of judges, or a jury) in order to render a decision about the origin of a speech voice recording whose identity is in question. Systems with lower EER or LLR suggest that they are more capable to generate instructive scores with higher validity and reliability. We explored several speaker adaptation methods and compared their speaker recognition performance, which aim to achieve a reliable system able to produce a measure of evidence in the form of a likelihood ratio (LR) score as the strength of evidence. The LR expresses the likelihood of the speech evidence under the two competing hypotheses (i.e. the prosecution hypothesis H_0 : the suspected speaker is the same as the source of the questioned recording versus the defense hypothesis H_1 : the suspected speaker is different from the source of the questioned recording [63]). The LR is the ratio between these two statements H_0 and H_1 .

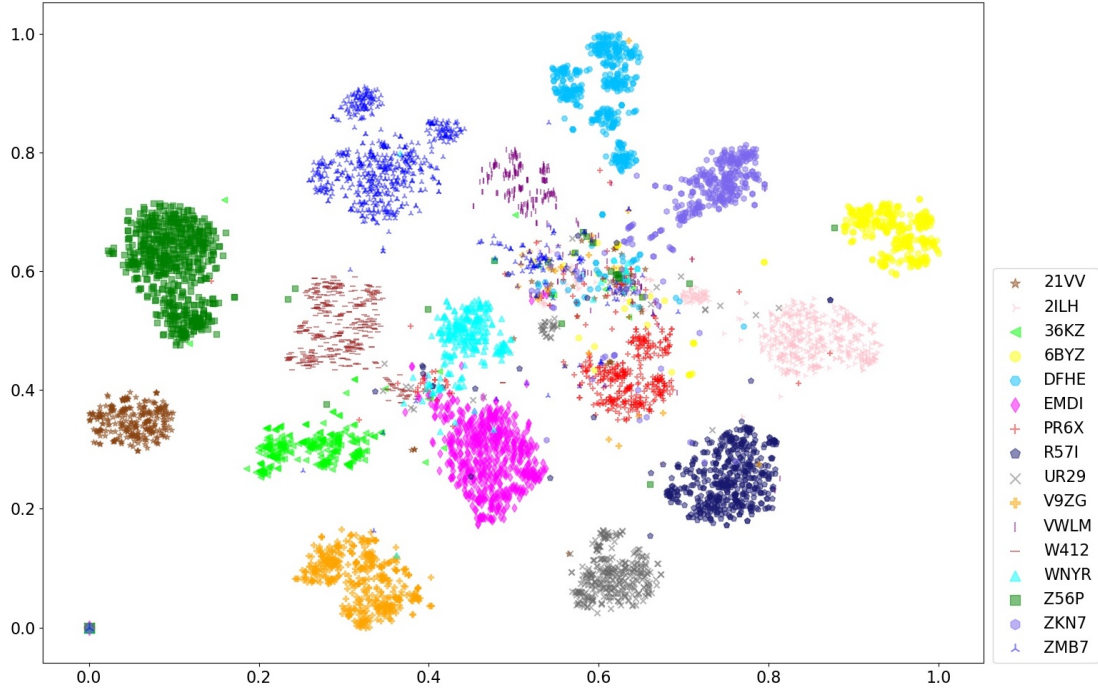
The x-Vector modeling approach with a PLDA backend can be applied for calculating LR. The goal of PLDA is to project data samples from the feature space to a latent space such that samples from the same class are modeled using the same distribution [45]. Given n utterance-level speaker embeddings $\{\mathbf{u}_i^p\}_{i=1}^n$ of speaker p in the latent space and one utterance-level speaker embedding \mathbf{u}^q of speaker q in the latent space, if we need to find whether they belong to same speaker or not, then we compute the likelihood ratio R based on two hypothesis H_0 and H_1 ,

$$R(\{\mathbf{u}_i^p\}_{i=1}^n, \mathbf{u}^q) = \frac{\text{likelihood}(H_0)}{\text{likelihood}(H_1)} = \frac{P(\{\mathbf{u}_i^p\}_{i=1}^n, \mathbf{u}^q)}{P(\{\mathbf{u}_i^p\}_{i=1}^n)P(\mathbf{u}^q)}, \quad (17)$$

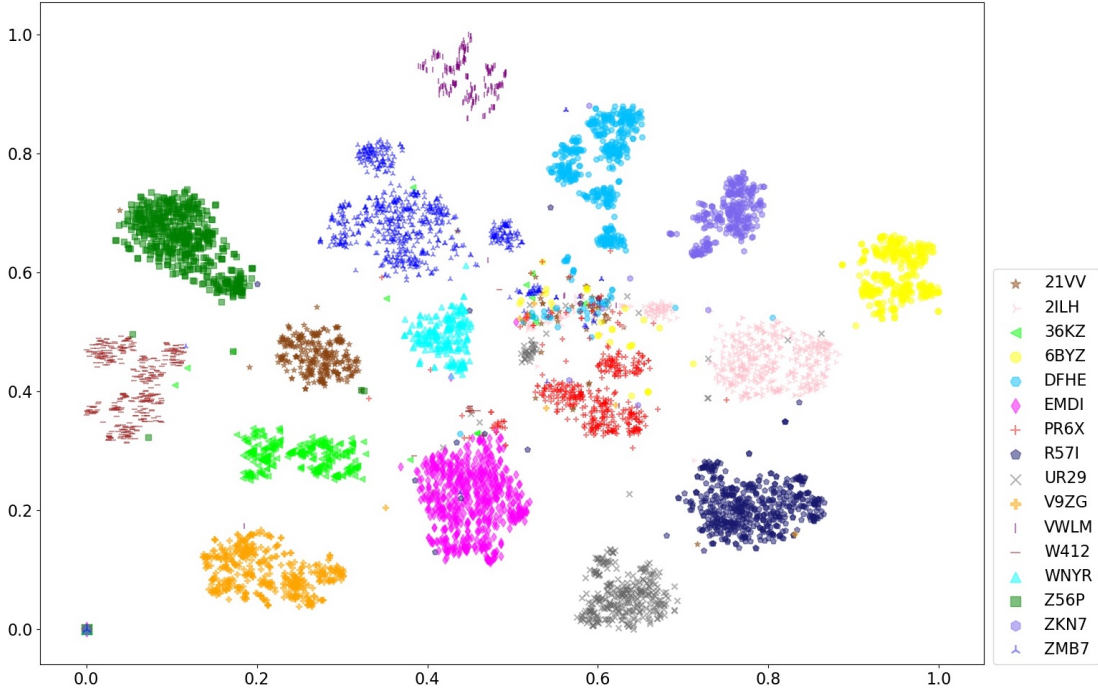
where

$$P(\{\mathbf{u}_i^p\}_{i=1}^n) = \int P(\mathbf{u}_1^p|\mathbf{v}) \dots P(\mathbf{u}_n^p|\mathbf{v})P(\mathbf{v})d\mathbf{v} \quad (18)$$

is the distribution of a set of examples, given that they belong to the same class, and \mathbf{v} represents the class centers in the latent space. The log likelihood ratio $\log(R)$ is known as PLDA scores. With larger $\log(R)$ values, there is stronger support for the H_0 hypothesis, and with smaller $\log(R)$ values, there is stronger support for the H_1 hypothesis. Here, we took 6 utterances for each speaker as enrollment data (here, $n = 6$ in Eq. (17)). Table IX gives an example of quantitative measurement in the form of PLDA score generated by the



(a) discrepancy-minimization



(b) moment-matching

Fig. 8: Visualization of LENA-field speaker embeddings by discrepancy-minimizing (a) and moment-matching (b) systems using t-SNE.

Moment-Matching system between the given speakers and a specific recording.

Speech data from these 16 speakers constitute the entire test set. The suspected speakers during LR calculation are called the relevant population in forensic speaker recognition. To avoid potential bias in the case proper, those speakers are often selected by a panel of listeners (e.g., police officers with linguistic background have no prior knowledge of a particular case) [64]. Since we already have the LR, it can be interpreted based on the odds form of Bayes' theorem, which

is represented as,

$$\frac{P(H_0|E)}{P(H_1|E)} = \frac{P(E|H_0)}{P(E|H_1)} \times \frac{P(H_0)}{P(H_1)}, \quad (19)$$

where E represents the observed speech evidence. This Bayes' theorem shows how the LR can be combined with prior knowledge concerning the case (knowledge unrelated to speech data) in order to arrive at posterior odds. Only the LR is the province provided by the speaker recognition system; the prior odds and posterior odds are the province of the court. The

TABLE IX
PLDA SCORES BETWEEN EACH SPEAKER AND A RECORDING

<i>speaker_id</i>	<i>record_id</i>	<i>PLDA_score</i>
21VV	21VV_LENA_field_100	7.585
21LH	21VV_LENA_field_100	-15.393
36KZ	21VV_LENA_field_100	0.145
6BYZ	21VV_LENA_field_100	-13.596
DFHE	21VV_LENA_field_100	-15.123
EMDI	21VV_LENA_field_100	-4.790
PR6X	21VV_LENA_field_100	-11.299
R57I	21VV_LENA_field_100	-9.353
UR29	21VV_LENA_field_100	-10.282
V9ZG	21VV_LENA_field_100	-0.368
VWLM	21VV_LENA_field_100	-1.464
W412	21VV_LENA_field_100	-5.462
WNYR	21VV_LENA_field_100	-5.688
Z56P	21VV_LENA_field_100	-4.420
ZKN7	21VV_LENA_field_100	-11.657
ZMB7	21VV_LENA_field_100	-5.653

forensic experts should only produce the LR in actual court cases and leave prior odds to the court or jury to interpret or assess. The judge or the jury in the court can use such a non-categorical opinion for their deliberations and decision.

IX. CONCLUSION

For forensic speaker recognition, addressing mismatch due to naturalistic field locations is a significant challenge. In general, fine-tuning is commonly employed for network model adaptation when a domain mismatch exists between train and test data. However, that approach usually considers only a single domain mismatch. In practical scenarios for forensic audio analysis, speech data are typically collected in multiple acoustic environments, which offer unique challenges to speaker recognition system development due to location uncertainty and diverse mismatch between reference and naturalistic field recordings. A speaker recognition system can deliver different speaker discrimination performance while evaluated on the dataset collected from multiple acoustic environments. In this study, we adopted a domain adversarial training (DAT) method with a gradient reversal layer to learn domain-invariant and speaker-discriminative representations. The DAT gives competitive results on each domain. Additionally, we formulated a discrepancy-minimizing solution to perform model adaptation for the purpose of improving speaker recognition performance across each potential field location with an overall smaller domain discrepancy. As demonstrated in our results, the solution improves speaker recognition system performance for each domain, which demonstrates that minimizing the domain discrepancy at both the frame-level and segment-level benefits system speaker discrimination. However, this improvement can still be unbalanced, as was shown with a higher EER result for the LENA-field set versus others. The LENA-field data was collected in locations entirely different from environments of the other three datasets. Accordingly, we proposed a moment-matching solution with an adversarial training schema for model adaptation to minimize domain discrepancy and simultaneously mitigate the impact of noise for LENA-field data

with the help of sound-booth captured audio. Consequently, the moment-matching system achieved the best speaker recognition results for each domain, with absolute EERs of 1.50%, 3.45%, 6.53%, 3.05% for the Clean, Far-field, LENA-field, and LENA-booth sets, respectively. Overall, the learned speaker representations through domain adversarial training (DAT), discrepancy-minimizing, and moment-matching solutions are less dependent on shifts in acoustic domains, which provides a solution to the challenging multi-source domain adaptation problem in forensic speaker recognition. Finally, we applied the most effective overall system for an independent simulative forensic case to show how the system solution can support the judge or jury in a court scenario to make a decision with a strength-of-evidence statement in the form of a likelihood ratio.

REFERENCES

- [1] J. H. L. Hansen and T. Hasan, "Speaker recognition by machines and humans: A tutorial review," *IEEE Signal processing magazine*, vol. 32, no. 6, pp. 74–99, 2015.
- [2] J. H. Wigmore, "New mode of identifying criminals," *Journal of Criminal Law and Criminology*, vol. 17, no. 2, pp. 155–156, 1926.
- [3] L. G. Kersta, "Voiceprint identification," *Nature*, vol. 196, no. 4861, pp. 1253–1257, 1962.
- [4] G. Pop, D. Draghicescu, and D. Burileanu, "A quality-aware forensic speaker recognition system," *Romanian Journal of Information Science and Technology*, vol. 17, no. 2, pp. 134–149, 2014.
- [5] L.-J. Boë, "Forensic voice identification in france," *Speech Communication*, vol. 31, no. 2-3, pp. 205–224, 2000.
- [6] Daniel-Ramos-Castro, "Forensic evaluation of the evidence using automatic speaker recognition systems," Ph.D. dissertation, Universidad autónoma de Madrid, 2007.
- [7] J.-F. Bonastre, F. Bimbot, L.-J. Boë, J. P. Campbell, D. A. Reynolds, and I. Magrin-Chagnolleau, "Person authentication by voice: A need for caution," in *Proc. ISCA INTERSPEECH-2003*, 2003, pp. 33–36.
- [8] E. Enzinger, "Implementation of forensic voice comparison within the new paradigm for the evaluation of forensic evidence," Ph.D. dissertation, University of New South Wales, Sydney, New South Wales, 2016.
- [9] N. Dehak, P. J. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, "Front-end factor analysis for speaker verification," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 4, pp. 788–798, 2010.
- [10] N. Dehak, P. A. Torres-Carrasquillo, D. Reynolds, and R. Dehak, "Language recognition via i-vectors and dimensionality reduction," in *Proc. ISCA INTERSPEECH-2011*, 2011, pp. 857–860.
- [11] M. Li and W. Liu, "Speaker verification and spoken language identification using a generalized i-vector framework with phonetic tokenizations and tandem features," in *Proc. ISCA INTERSPEECH-2014*, 2014, pp. 1120–1124.
- [12] C. Zhang, K. Koishida, and J. H. Hansen, "Text-independent speaker verification based on triplet convolutional neural network embeddings," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 9, pp. 1633–1644, 2018.
- [13] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, "X-vectors: Robust dnn embeddings for speaker recognition," in *Proc. 2018 IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2018, pp. 5329–5333.
- [14] D. Snyder, P. Ghahremani, D. Povey, D. Garcia-Romero, Y. Carmiel, and S. Khudanpur, "Deep neural network-based speaker embeddings for end-to-end speaker verification," in *Proc. 2016 IEEE Spoken Language Technology Workshop (SLT)*, IEEE, 2016, pp. 165–170.
- [15] C. Li, X. Ma, B. Jiang, X. Li, X. Zhang, X. Liu, Y. Cao, A. Kannan, and Z. Zhu, "Deep speaker: an end-to-end neural speaker embedding system," *arXiv preprint arXiv:1705.02304*, vol. 650, 2017.
- [16] Y. Lei, N. Scheffer, L. Ferrer, and M. McLaren, "A novel scheme for speaker recognition using a phonetically-aware deep neural network," in *Proc. 2014 IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2014, pp. 1695–1699.
- [17] W. Cai, J. Chen, and M. Li, "Exploring the encoding layer and loss function in end-to-end speaker and language recognition system," in *Proc. Odyssey*, pp. 74–81, 2018.

- [18] M. Kumar, T. Jin-Park, S. Bishop, and S. Narayanan, "Designing neural speaker embeddings with meta learning," *arXiv preprint arXiv:2007.16196*, 2020.
- [19] W. Xia, C. Zhang, C. Weng, M. Yu, and D. Yu, "Self-supervised text-independent speaker verification using prototypical momentum contrastive learning," *arXiv preprint arXiv:2012.07178*, 2020.
- [20] J. H. Hansen and V. Varadarajan, "Analysis and compensation of lombard speech across noise type and levels with application to in-set/out-of-set speaker recognition," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, no. 2, pp. 366–378, 2009.
- [21] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Deep transfer learning with joint adaptation networks," in *Proc. Int. Conf. on machine learning*, 2017, pp. 2208–2217.
- [22] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell, "Deep domain confusion: Maximizing for domain invariance," *arXiv preprint arXiv:1412.3474*, 2014.
- [23] M. Long, Y. Cao, J. Wang, and M. Jordan, "Learning transferable features with deep adaptation networks," in *Proc. Int. Conf. on Machine Learning*, 2015, pp. 97–105.
- [24] W. Zellinger, T. Grubinger, E. Lughofer, T. Natschlager, and S. Saminger-Platz, "Central moment discrepancy (cmd) for domain-invariant representation learning," *arXiv preprint arXiv:1702.08811*, 2017.
- [25] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2017, pp. 7167–7176.
- [26] Y. Ganin and V. Lempitsky, "Unsupervised domain adaptation by backpropagation," in *Proc. Int. Conf. on Machine Learning*, 2015, pp. 1180–1189.
- [27] K. Saito, K. Watanabe, Y. Ushiku, and T. Harada, "Maximum classifier discrepancy for unsupervised domain adaptation," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2018, pp. 3723–3732.
- [28] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proc. IEEE Int. Conf. on Computer Vision*, 2017, pp. 2223–2232.
- [29] M.-Y. Liu, T. Breuel, and J. Kautz, "Unsupervised image-to-image translation networks," in *Proc. Advances in Neural Information Processing Systems*, 2017, pp. 700–708.
- [30] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, and T. Darrell, "Cycada: Cycle-consistent adversarial domain adaptation," in *Proc. Int. Conf. on Machine Learning*, 2018, pp. 1989–1998.
- [31] G. Bhattacharya, J. Monteiro, J. Alam, and P. Kenny, "Generative adversarial speaker embedding networks for domain robust end-to-end speaker verification," in *Proc. 2019 IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6226–6230.
- [32] W. Xia, J. Huang, and J. H. L. Hansen, "Cross-lingual text-independent speaker verification using unsupervised adversarial discriminative domain adaptation," in *Proc. 2019 IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 5816–5820.
- [33] W. Lin, M.-M. Mak, N. Li, D. Su, and D. Yu, "Multi-level deep neural network adaptation for speaker verification using mmd and consistency regularization," in *Proc. 2020 IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6839–6843.
- [34] W. Lin, M.-M. Mak, N. Li, D. Su, and D. Yu, "A framework for adapting dnn speaker embedding across languages," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 2810–2822, 2020.
- [35] S. S. Sarfjoo, S. Marcel *et al.*, "Domain adaptation and investigation of robustness of dnn-based embeddings for text-independent speaker verification using dilated residual networks," IDIAP, Tech. Rep., 2019.
- [36] K. Vesely, S. Watanabe, K. Žmolíková, M. Karafiát, L. Burget, and J. H. Černocký, "Sequence summarizing neural network for speaker adaptation," in *Proc. 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 5315–5319.
- [37] Z. Wang, W. Xia, and J. H. L. Hansen, "Cross-domain adaptation with discrepancy minimization for text-independent forensic speaker verification," *Proc. ISCA INTERSPEECH-2020*, pp. 2257–2261, 2020.
- [38] Y. Zhu, F. Zhuang, J. Wang, J. Chen, Z. Shi, W. Wu, and Q. He, "Multi-representation adaptation network for cross-domain image classification," *Neural Networks*, vol. 119, pp. 214–221, 2019.
- [39] Y. Zhang, B. Deng, H. Tang, L. Zhang, and K. Jia, "Unsupervised multi-class domain adaptation: Theory, algorithms, and practice," *IEEE transactions on pattern analysis and machine intelligence*, vol. 14, no. 8, 2020.
- [40] R. Xu, Z. Chen, W. Zuo, J. Yan, and L. Lin, "Deep cocktail network: Multi-source unsupervised domain adaptation with category shift," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3964–3973.
- [41] X. Peng, Q. Bai, X. Xia, Z. Huang, K. Saenko, and B. Wang, "Moment matching for multi-source domain adaptation," in *Proc. IEEE International Conference on Computer Vision*, 2019, pp. 1406–1415.
- [42] Y. Zhu, F. Zhuang, and D. Wang, "Aligning domain-specific distribution and classifier for cross-domain classification from multiple sources," in *Proc. AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 5989–5996.
- [43] V. Peddinti, D. Povey, and S. Khudanpur, "A time delay neural network architecture for efficient modeling of long temporal contexts," in *Proc. ISCA INTERSPEECH-2015*, 2015, pp. 3214–3218.
- [44] D. Snyder, D. Garcia-Romero, D. Povey, and S. Khudanpur, "Deep neural network embeddings for text-independent speaker verification," in *Proc. ISCA INTERSPEECH-2017*, 2017, pp. 999–1003.
- [45] S. Ioffe, "Probabilistic linear discriminant analysis," in *Proc. European Conf. on Computer Vision*. Springer, 2006, pp. 531–542.
- [46] D. E. Sturim and D. A. Reynolds, "Speaker adaptive cohort selection for tnrm in text-independent speaker verification," in *Proc. 2005 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, vol. 1, 2005, pp. 741–744.
- [47] J. P. Campbell, W. Shen, W. M. Campbell, R. Schwartz, J.-F. Bonastre, and D. Matrouf, "Forensic speaker recognition," *IEEE Signal Processing Magazine*, vol. 26, no. 2, pp. 95–103, 2009.
- [48] J. S. Chung, A. Nagrani, and A. Zisserman, "Voxceleb2: Deep speaker recognition," pp. 999–1003, 2018.
- [49] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," *The journal of machine learning research*, vol. 17, no. 1, pp. 2096–2030, 2016.
- [50] A. Gretton, K. Borgwardt, M. Rasch, B. Schölkopf, and A. Smola, "A kernel method for the two-sample-problem," *Proc. Advances in Neural Information Processing Systems*, vol. 19, pp. 513–520, 2006.
- [51] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola, "A kernel two-sample test," *The Journal of Machine Learning Research*, vol. 13, no. 1, pp. 723–773, 2012.
- [52] Y. Li, K. Swersky, and R. Zemel, "Generative moment matching networks," in *Proc. Int. Conf. on Machine Learning*, 2015, pp. 1718–1727.
- [53] D. Snyder, G. Chen, and D. Povey, "Musan: A music, speech, and noise corpus," *arXiv preprint arXiv:1510.08484*, 2015.
- [54] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.
- [55] S. J. Prince and J. H. Elder, "Probabilistic linear discriminant analysis for inferences about identity," in *Proc. 2007 IEEE 11th Int. Conf. on Computer Vision*. IEEE, 2007, pp. 1–8.
- [56] A. Martin, G. Doddington, T. Kamm, M. Ordowski, and M. A. Przybicki, "The det curve in assessment of detection task performance," in *Proc. ESCAPEUROSPPEECH*, 1997, pp. 1895–1898.
- [57] N. Brümmer and J. Du Preez, "Application-independent evaluation of speaker detection," *Computer Speech & Language*, vol. 20, no. 2-3, pp. 230–275, 2006.
- [58] F. Kelly and J. H. Hansen, "Score-aging calibration for speaker verification," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 12, pp. 2414–2424, 2016.
- [59] D. A. Van Leeuwen and N. Brümmer, "An introduction to application-independent evaluation of speaker recognition systems," in *Speaker classification I*. Springer, 2007, pp. 330–353.
- [60] S. Pigeon, P. Druyts, and P. Verlinde, "Applying logistic regression to the fusion of the nist'99 1-speaker submissions," *Digital Signal Processing*, vol. 10, no. 1-3, pp. 237–248, 2000.
- [61] A. Alexander, O. Forth, J. Nash, and N. Yager, "Zooplots for speaker recognition with tall and fat animals," *Proc. Intl. Assoc. for Forensic Phonetics and Acoustics (IAFPA)*, [Online] <http://www.oxfordwaveresearch.com/papers/Alexander-Forth-Nash-Yager-IAFPA-2014-Abstract.pdf>, accessed, pp. 11–29, 2017.
- [62] L. v. d. Maaten and G. Hinton, "Visualizing data using t-sne," *Journal of Machine Learning Research*, vol. 9, no. Nov, pp. 2579–2605, 2008.
- [63] A. Drygajlo, "Automatic speaker recognition for forensic case assessment and interpretation," in *Forensic Speaker Recognition*. Springer, 2012, pp. 21–39.
- [64] G. S. Morrison, F. Ochoa, and T. Thiruvanan, "Database selection for forensic voice comparison," in *Odyssey 2012-The Speaker and Language Recognition Workshop*, 2012.



Zhenyu Wang received the B.S. degree majoring in digital media technology from Hangzhou Dianzi University, Hangzhou, China in 2015. He received the M.S. degree in engineering computer system architecture from Beijing Language and Culture University, Beijing, China in 2019. He started working toward the Ph.D. degree in computer engineering at the University of Texas at Dallas (UTD), Richardson, TX, USA, in 2019. He works with Professor John H.L. Hansen at the Center for Robust Speech Systems. Since the same year, he has been a Graduate

Research Assistant with the Center for Robust Speech Systems (CRSS), UTD. His research interests include mispronunciation verification, computer-assisted language learning, forensic audio analysis and model adaptation for open-set speaker recognition system, representation learning used for acoustic modeling. He has authored around five journal and conference papers in the field of speech processing and language technology.



John H. L. Hansen John H. L. Hansen (S'81–M'82–SM'93–F'07) received the B.S.E.E. degree from the College of Engineering, Rutgers University, New Brunswick, NJ, USA, in 1982, the M.S. and Ph.D. degrees in electrical engineering from Georgia Institute of Technology, Atlanta, GA, USA, in 1983 and 1988, respectively. He joined Erik Jonsson School of Engineering and Computer Science, University of Texas at Dallas (UTDallas), in 2005, where he currently serves as a Jonsson School Associate Dean for Research, as well as a Professor

of Electrical and Computer Engineering, the Distinguished University Chair in Telecommunications Engineering, and a joint appointment as a Professor in the School of Behavioral and Brain Sciences (Speech & Hearing). He previously served as a UTDallas Department Head of Electrical Engineering from August 2005 to December 2012, overseeing a +4x increase in research expenditures (4.5 to 22.3 M) with a 20% increase in enrollment along with hiring 18 additional T/TT faculty, growing UTDallas to the eighth largest EE program from ASEE rankings in terms of degrees awarded. At UTDallas, he established the Center for Robust Speech Systems (CRSS). Previously, he served as a Department Chairman and Professor of Speech, Language and Hearing Sciences, and a Professor in Electrical and Computer Engineering, University of Colorado-Boulder (1998–2005), where he co-founded and served as an Associate Director of the Center for Spoken Language Research. In 1988, he established the Robust Speech Processing Laboratory and continued to direct research activities in CRSS at UTDallas. He is the author/coauthor of 800+ journal and conference papers including 13 textbooks in the field of speech processing and language technology, signal processing for vehicle systems, coauthor of textbook: *Discrete-Time Processing of Speech Signals* (IEEE Press, 2000), co-editor of *DSP for In-Vehicle and Mobile Systems* (Springer, 2004), *Advances for In-Vehicle and Mobile Systems: Challenges for International Standards* (Springer, 2006), *In-Vehicle Corpus and Signal Processing for Driver Behavior* (Springer, 2008), and lead author of the report *The Impact of Speech Under Stress on Military Speech Technology*, (NATO RTOTR-10, 2000). He has supervised 95 Ph.D./M.S. thesis candidates (54 Ph.D., 41 M.S./M.A.). His research interests include the areas of digital speech processing, analysis and modeling of speech and speaker traits, speech enhancement, feature estimation in noise, signal processing for hearing impaired/cochlear implants, robust speech recognition with emphasis on machine learning and knowledge extraction, and in-vehicle interactive systems for hands-free human-computer interaction. Dr. Hansen received the honorary degree Doctor Technicus Honoris Causa from Aalborg University (Aalborg, DK) in April 2016, in recognition of his contributions to speech signal processing and speech/language/hearing sciences. He was recognized as an IEEE Fellow (2007) for contributions in “Robust Speech Recognition in Stress and Noise,” International Speech Communication Association (ISCA) Fellow (2010) for contributions on research for speech processing of signals under adverse conditions, and received The Acoustical Society of America 25 Year Award (2010) in recognition of his service, contributions, and membership to the Acoustical Society of America. He previously served as ISCA President (2017–2021) and Vice-president (2015–2017) and currently serves as tenure and member of the ISCA Board, having previously served as the Vice-President (2015–2017). He also is serving as a Vice-Chair on U.S. Office of Scientific Advisory Committees (OSAC) for OSAC-Speaker in the voice forensics domain (2015–2017). Previously, he served as an IEEE Technical Committee (TC) Chair and member of the IEEE Signal Processing Society: Speech-Language Processing Technical Committee (SLTC) (2005–2008; 2010–2014; elected IEEE SLTC Chairman for 2011–2013, Past-Chair for 2014), and elected as an ISCA Distinguished Lecturer (2011–2012). He has served as the member of the IEEE Signal Processing Society Educational Technical Committee (2005–2008; 2008–2010); Technical Advisor to the U.S. Delegate for NATO (IST/TG-01); IEEE Signal Processing Society Distinguished Lecturer (2005–2006), Associate Editor for IEEE TRANSACTION SPEECH AND AUDIO PROCESSING (1992–1999), Associate Editor for IEEE SIGNAL PROCESSING LETTERS (1998–2000), Editorial Board Member for IEEE SIGNAL PROCESSING MAGAZINE (2001–2003); and Guest Editor (October 1994) for special issue on Robust Speech Recognition for IEEE TRANSACTION SPEECH AND AUDIO PROCESSING. He is serving as an Associate Editor for JASA, and served on Speech Communications Technical Committee for Acoustical Society of America (2000–2003). He was the recipient of The 2005 University of Colorado Teacher Recognition Award as voted on by the student body. He organized and served as the General Chair for ISCA INTERSPEECH-2002, September 16–20, 2002, Co-Organizer, and Technical Program Chair for IEEE ICASSP-2010, Dallas, TX, USA, March 15–19, 2010, and Co-Chair and Organizer for IEEE SLT-2014, December 7–10, 2014 in Lake Tahoe, NV, USA.