
SpeakerGAN: Recognizing Speakers in New Languages with Generative Adversarial Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Verifying a person’s identity based on their voice, is a challenging, real-world
2 problem in biometric security. A crucial requirement of such speaker verification
3 systems is to be domain robust. Performance should not degrade even if speakers
4 are talking in languages not seen during training. To this end, we propose to use
5 Generative Adversarial Networks to adapt a speaker embedding model to new
6 languages using a small amount of unlabelled data. Our model is optimized end-to-
7 end, and verification can be performed using simple cosine scoring. Additionally,
8 we propose a novel objective for training the generator, that yields further im-
9 provements. We show that the proposed unsupervised adversarial adaptation leads
10 to verification performance that is competitive with state-of-the-art verification
11 systems. In an attempt to better understand the performance of our models, we
12 quantitatively measure the degree of invariance induced by our proposed methods
13 using Maximum Mean Discrepancy and Fréchet distances.

14 1 Introduction

15 Text-Independent Speaker Verification remains a challenging problem in the domain of biometric
16 security. Armed with the machinery of deep learning, verification systems can now be deployed
17 in the wild, and are still capable of delivering robust performance. In the verification community,
18 situations wherein the test data is significantly different from the data available during system training
19 are referred to as - In the Wild. For instance, the NIST-SRE 2016 evaluation data contains Cantonese
20 and Tagalog speakers (in-domain, target data), while most of the speakers in our training set are
21 talking in English (out-of-domain, source data). This distribution shift or mismatch between training
22 and test data is an obstacle in several areas of pattern recognition and machine learning [1], and leads
23 to a degradation in system performance. The development biometric verification system that perform
24 reliably in such conditions is critical for this technology be used safely and securely on a day-to-day
25 basis.

26 Deep neural networks (DNN) have revolutionized several areas of speech processing, and as such,
27 are ideal candidates for learning discriminative speaker representations or embeddings [11, 4, 13, 2].
28 Indeed, neural speaker embeddings have surpassed the performance of i-vectors [11, 3], especially
29 on real world, in the wild data [8, 6]. Arguably the most popular approach for learning speaker
30 embeddings is to optimize the parameters of a DNN by minimizing the cross-entropy loss over
31 speakers in the training data. Cross-entropy is natural choice for identifying speakers, however it
32 does not directly address the verification task. As a consequence of not being optimized ‘end-to-
33 end’, the performance of cross-entropy speaker embeddings (X-vectors) is heavily dependent on a
34 powerful classifier to perform verification. This dependence on a classifier motivates the research and
35 development of end-to-end systems. We also believe that such systems can also benefit in downstream
36 tasks that make use of speaker embeddings, such as speech recognition and synthesis.

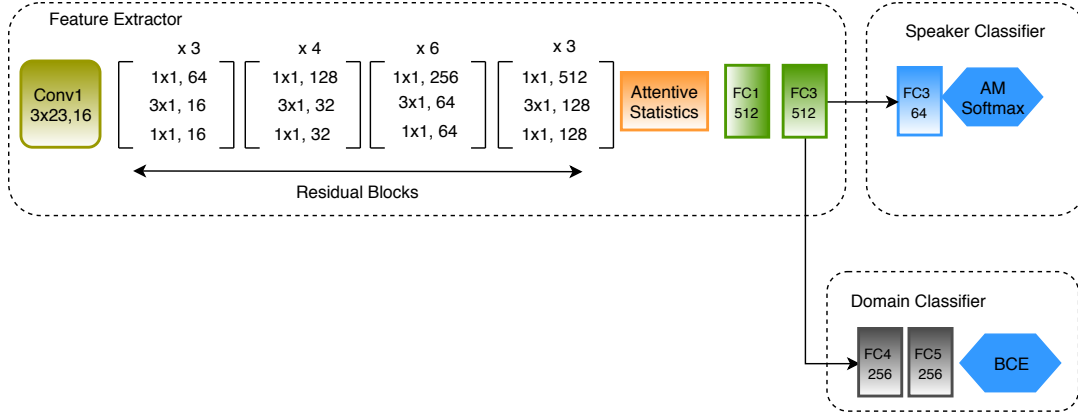


Figure 1: Domain Adversarial Neural Speaker Embedding Model

Verifying a speaker’s identity is a challenging problem. Modern speaker verification datasets like NIST-SRE 2016, add to this challenge by introducing a mismatch between the distributions of the training and test data. This phenomena is referred to as domain or covariate shift. In the case of NIST-SRE 2016, the test data consists of Cantonese and Tagalog speakers, whereas the vast majority of training speakers are talking in English. NIST also provide a small amount of unlabelled, *in-domain*, *target* data, that can be used to compensate for the domain shift. Most the domain adaptation techniques that have been proposed for speaker verification have been proposed on top of i-vectors or x-vectors.

In this we present a unified framework for directly learning domain-robust speaker embeddings using Generative Adversarial Networks (GAN). We draw inspiration from research in computer vision, where GAN based domain adaptation methods have been shown to be more powerful than the gradient reversal framework [12, 9, 10]. Both methods cast domain adaptation/invariance as an adversarial game - generate features or embeddings such that a discriminator cannot tell if they come from the source or target domain. Unlike traditional GANs that work in high-dimensional spaces (e.g. natural images, speech), domain adaptation GANs operate in low-dimensional embedding space. Keeping these constraints in mind, we propose a novel objective for updating the generator network, which we find to work better than the conventional generative loss.

The nature of the adversarial game makes training GAN models challenging. We found that a simple way to stabilize the training of our models was to make the GANs conditional. Specifically we propose to use a modified version of the Auxiliary Classifier GAN (AuxGAN)[]. We show that this addition makes our model robust and we are able to use the same set of hyper-parameters to train several GAN variants within our framework. In our experiments we show that all of the models outperform the DANSE model, with some delivering comparable performance to a state-of-the-art x-vector system. Furthermore, we are able to fuse different GAN models using simple score averaging to achieve state-of-the-art verification performance.

2 Models

2.1 Feature Extractor (Generator)

The first step for learning discriminative speaker embeddings is to learn a mapping $F(X_s) \rightarrow \mathbf{f}$, $\mathbf{f} \in R^D$ from a sequence of speech frames from speaker s to a D-dimensional feature vector \mathbf{f} . $F(X)$ can be implemented using a variety of neural network architectures. We design our feature extractor using a residual network structure. We choose to model speech using 1-dimensional convolutional filters, owing to the fact that speech is translation invariant along the time-axis only. Following the residual blocks we use a combination of self-attention and dense layers in order to represent input audio of arbitrary size by a fixed-size vector, \mathbf{f} . Unlike traditional approaches, our proposed feature extractor is updated with an adversarial loss in addition to the standard task loss.

72 2.2 Self-Attentive Speaker Statistics

73 Self-Attention models are an active area of research in the speaker verification community. Intuitively,
 74 such models allow the network to focus on fragments of speech that are more speaker discriminative.
 75 The attention layers computes a scalar weight corresponding to each time-step t :

$$e_t = \mathbf{v}^T f(\mathbf{W}h_t + \mathbf{b}) + k \quad (1)$$

76 These weights are then normalized, $\alpha_t = \text{softmax}(e_t)$, to give them a probabilistic interpretation.
 77 We use the attention model proposed in [15], which extends attention to the mean as well as standard
 78 deviation:

$$\hat{\mu} = \sum_t^T \alpha_t \mathbf{h}_t \quad (2)$$

$$\hat{\sigma} = \sum_t^T \alpha_t \mathbf{h}_t \odot \mathbf{h}_t - \hat{\mu} \odot \hat{\mu} \quad (3)$$

79 In this work we apply the use of self attention to convolutional feature maps, as indicated in Fig. 1.
 80 The last residual block outputs a tensor of size $nB \times nF \times T$, where nB is the batch size, nF is the
 81 number of filters and T is time. The input to the attention layer, h_t , is a nF dimensional vector.

82 By using a self-attention model, we also equip our network with a more robust framework for
 83 processing inputs of arbitrary size than simple global averaging. This allows us simply forward
 84 propagate a recording through the network in order to extract speaker embeddings.

85 2.3 Classifier

86 The classifier block, $C(\mathbf{f}, \theta_y)$, is arguably the key component of the model, as it is responsible for
 87 learning speaker discriminative features. Recently, angular margin loss functions have been proposed
 88 as an alternative to contrastive loss functions for verification tasks [5, 14]. The Additive Margin
 89 softmax (AM-softmax) loss function is one such algorithm with an intuitive interpretation. The loss
 90 computes similarity between classes using cosine, and forces the similarity of the correct class to be
 91 greater than that of incorrect classes by a margin m .

$$\begin{aligned} L_{AMS} &= -\frac{1}{n} \sum_{i=1}^n \log \frac{e^{s \cdot (\cos \theta_{y_i} - m)}}{e^{\cos \theta_{s \cdot (y_i - m)}} + \sum_{j \neq y_i} e^{s \cdot (\cos \theta_j)}} \\ &= -\frac{1}{n} \sum_{i=1}^n \log \frac{e^{s \cdot (W^T \mathbf{f}_i - m)}}{e^{s \cdot (W^T \mathbf{f}_i - m)} + \sum_{j \neq y_i} e^{s \cdot (W^T \mathbf{f}_j)}} \end{aligned} \quad (4)$$

92 Where W^T and \mathbf{f}_i are the normalized weight vector and speaker embedding respectively. The
 93 AM-softmax loss also adds a scale parameter s , which helps the model converge faster. We select
 94 $m = 0.6$ and $s = 30$ for all our experiments.

95 2.4 Domain Discriminator

96 The domain discriminator $D(\cdot)$ is tasked with determining if embeddings come from the source
 97 or target domains, and is arguably the most important component of the model. In order to learn
 98 domain invariant features, we engage the domain discriminator in an adversarial game with the feaure
 99 extractor $E(\cdot)$. The domain discriminator consists of two fully connected layers followed by the
 100 output layer.

101 3 Domain Adversarial Speaker Embeddings

102 A key requirement for learning speaker embeddings that are domain invariant is to find a balance
 103 between the task loss and the adversarial loss. The objective to learn a feature space wherein

104 embeddings are speaker discriminative irrespective of the domain. Key to achieving this is the domain
 105 discriminator D , which is trained using the Binary Cross-Entropy loss (BCE).

$$\mathcal{L}_{adv_D}(\mathbf{X}_s, \mathbf{X}_t, E) = -E_{x_s \sim X_s}[\log(D(E(x_s)))] - E_{x_t \sim X_t}[\log(1 - D(E(x_t)))] \quad (5)$$

106 Where $\mathbf{X}_s, \mathbf{X}_t$ represent source and target data respectively. $E(\cdot)$ is the feature extractor/generator.
 107 The adversarial game between $D(\cdot)$ and $E(\cdot)$ is given by:

$$\begin{aligned} \min_D \mathcal{L}_{adv_D}(\mathbf{X}_s, \mathbf{X}_t, E) \\ \min_E \mathcal{L}_{adv_E}(\mathbf{X}_s, \mathbf{X}_t, D) \end{aligned} \quad (6)$$

108 Equation (3) represents the most general form of the GAN game, and can be used to represent
 109 different adversarial frameworks depending on the choice of \mathcal{L}_{adv_E}

110 **Gradient Reversal:** We obtain the gradient reversal framework by setting $\mathcal{L}_{adv_E} = -\mathcal{L}_{adv_D}$.
 111 Gradient reversal optimizes the true minmax objective of the adversarial game. However, this
 112 objective can become problematic, since the discriminator converges early during training and leads
 113 to vanishing gradients.

114 **GAN:** Rather than directly using the minimax loss, the standard way to train the generator is using
 115 the inverted label loss. The generator objective is given by:

$$\mathcal{L}_{adv_E}(\mathbf{X}_s, \mathbf{X}_t, D) = -\mathbb{E}_{x_s \sim X_s}[\log(D(E(x_t)))] \quad (7)$$

116 This splits the optimization into two independent objectives, one for the generator and one for the
 117 discriminator. This loss has the same fixed-point properties as the minimax loss while providing
 118 stronger gradients to target mappings [12].

119 3.1 Updating the Generator with Source Embeddings

120 In a typical GAN setting, the generator is trained only using fake data (with inverted labels). This
 121 structure is also maintained in several adversarial domain adaptation algorithms []. However, in the
 122 context of this work we believe that updating the generator using *both* source and target data can be
 123 beneficial. In this case, the generator loss simply inverts the discriminator loss of eq. (1):

$$\begin{aligned} \mathcal{L}_{adv_E}(\mathbf{X}_s, \mathbf{X}_t, D) = \\ -E_{x_s \sim X_s}[\log(D(E(x_t)))] \\ -E_{x_t \sim X_t}[\log(1 - D(E(x_s)))] \end{aligned} \quad (8)$$

124 When using the proposed objective for training the generator, we are optimizing the true minimax
 125 loss like in the gradient reversal approach. Unfortunately, we found that optimizing this loss becomes
 126 unstable early during training. We found a simple approach to stabilize training for this model was to
 127 augment the discriminator with an auxiliary loss function. This addition also makes training more
 128 stable when optimizing the standard generator objective.

129 3.2 Auxiliary Classifier GAN

130 The Auxiliary Classifier GAN (AuxGAN) model augments the standard GAN framework with an
 131 auxiliary loss to perform conditional image generation [7]. This approach aims to predict side
 132 information (such as class labels), as opposed to feed the same information to the generator and
 133 discriminator. In the context of this work, our goal is to use the prediction loss for regularization and
 134 representation learning.

$$\begin{aligned} \min_D \mathcal{L}_{adv_D}(\mathbf{X}_s, \mathbf{X}_t, E) + \mathcal{L}_{Aux}(\mathbf{X}_s, Y_s) \\ \min_E \mathcal{L}_{adv_E}(\mathbf{X}_s, \mathbf{X}_t, D) + \mathcal{L}_{Aux}(\mathbf{X}_s, Y_s) \end{aligned} \quad (9)$$

Eq. (6) is a modified version of the AuxGAN objective. In particular, the original formulation also uses the auxiliary loss to train the generator as well (with fake data being assigned its own unique label). We found that the auxiliary loss was crucial for stabilizing $\mathcal{L}_{adv_E}(\mathbf{X}_s, \mathbf{X}_t, D)$ when using the formulation in eq. (6). In our experiments we found that the AuxGAN setup stabilizes model training even when we use eq. (3) as the generator objective, and leads to slightly better verification performance. In this setting only the discriminator is trained with the auxiliary loss.

3.3 GAN Variants

Since their introduction, GANs have been one of the most researched topics in the deep learning community. Several variations of the original formulation have been proposed, each with different generative characteristics and stability issues. In this work we explore three GAN variants in addition to the standard GAN - Least-Squares GAN [], Auxiliary Classifier GAN and Relativistic GAN []. These models differ in the structure of the discriminator network. We show that each variant transforms the feature space in different way, will all the model showing mostly similar performance. Additionally we see that by fusing the performance of all GAN variants together through score averaging we achieve the best overall performance.

4 Experimental Setup

Training Data (Source): We used audio from previous NIST-SRE evaluations (2004-2010) and Switchboard Cellular audio for training the proposed DANSE model as well as the x-vector and i-vector baseline systems. We also augment our data with noise and reverberation, as in []. We add 128k noisy copies to the clean speech, ending up with 220k recordings in our training set. For DANSE model training we filter out speakers with less than 5 recordings, ending up with approximately 6000 speakers, whereas the x-vector and i-vector systems were trained using the Kaldi recipe. We note that the vast majority of our training data consists of English speakers, and is recorded over telephone/cellular channels.

Model: In order to make a fair comparison, we use an identical network to the DANSE model. The *Embedding function/Generator*, $E(\cdot)$, consists of a 3×23 input convolutional layer, 4 residual blocks [3,4,6,3], an attentive statistics layer and two fully connected layers (512,512). The *classifier*, $C(\cdot)$, module consists of a fully connected layer (64) and the AM-softmax output layer. The former is the final domain invariant speaker embedding extracted during evaluation. Finally, the *domain discriminator* module consists of two fully connected layers (256,256) and a binary cross-entropy output layer. Exponential Linear Units (ELU) are used as non-linear activations for all layers of the network. Batch Normalization is used on all layers except the attentive statistics layer. We refer the reader to [] for a detailed description of the model.

Optimization: We start by pre-training the Embedding function using standard cross-entropy training. Pre-training is carried out using the RMSprop optimizer with a learning rate (lr) of 0.001. For training GAN based speaker embedding models we use different optimizers for training the three networks (Embedding function, Classifier, Discriminator). The classifier is optimized using RMSprop with $lr=0.003$, while the domain classifier and feature extractor are trained using SGD with $lr=0.001$. We were able to train all our GAN models using the same set of hyper-parameters. We used performance on held out validation set to determine when to stop training.

Data Sampling: We use an extremely simple approach for sampling data during training. We sample random chunks of audio (3-8 seconds) from each recording in the training set. We sample each recording 10 times to define an epoch. For each mini-batch of source data, we randomly sample (with repetition) a mini-batch from the unlabelled adaptation data for GAN training.

Speaker Verification: At test time we discard the domain discriminator branch of the model, as it is not needed for extracting embeddings. Extraction is done by performing a forward pass on the full recording, and using the 64-dimensional *fc3* layer as our speaker embeddings. Verification trials are scored using cosine distance. Verification performance is reported in terms of Equal Error Rate (EER).

5 Results

NIST-SRE 2016: Unlike previous years, The 2016 edition of the NIST-SRE introduced a challenging new dataset containing Cantonese and Tagalog speakers. We use the Kaldi recipes for our baseline i-vector and x-vector systems. We note that the x-vector baseline may be considered as state-of-the-art performance on this dataset.

Adaptation Data (Target): 2272 unlabelled, target data recordings are provided to adapt verification systems.

Table 1: Baseline Systems

<i>Model</i>	<i>Classifier</i>	<i>Cantonese</i>	<i>Tagalog</i>	<i>Pooled</i>
i-vector	PLDA	B	C	D
x-vector	COSINE	36.44	41.07	38.69
x-vector	LDA/PLDA	7.03	15.41	11.15
x-vector	PLDA	18.46	7.99	13.32
DASE	COSINE	17.87	8.84	13.36

Table 2: Performance of Different GAN systems in terms of EER(%). **SGAN:** standard, **AuxGAN:** auxiliary classifier, **LSGAN:** least squares, **RelGAN:** reletavistic, **FuseGAN:** score averaging

<i>Model</i>	<i>Classifier</i>	<i>Cantonese</i>	<i>Tagalog</i>	<i>Pooled</i>
SGAN	COSINE	8.32	17.51	D
AuxGAN	COSINE	7.88	16.10	11.93
LSGAN	COSINE	7.92	15.63	11.74
RelGAN	COSINE	8.01	16.22	13.32
FuseGAN	COSINE	6.93	14.84	10.88

Table 1. compares the performance of the different speaker representations on the NIST-SRE16 task. The x-vector model produces the best results, however requires LDA based dimensionality reduction and the PLDA classifier. The DASE model produces competitive results, showing similar performance to the x-vector/PLDA system that does not use LDA. Our model also performs slightly better than the i-vector/PLDA system. We note that the PLDA classifier also requires significant tuning and data augmentation. Furthermore, different PLDA implementations can lead to significantly different (worse) verification performance.

6 Analysis

One particularly interesting result from our experiments is the improvement we see through a simple score averaging procedure. Our hypothesis is that the different discriminator objectives encourage the generator to cover different modes of the target data distribution. This finding is consistent with GAN approaches that train multiple discriminators [], although we do not train them simultaneously. In Fig. 2 we visualize the embedding spaces learned by our models using t-SNE []. While Gradient Reversal primarily appears to rotate the feature space, the transformations induced by the GAN models is more pronounced. Crucially, we see that that the source domain speaker clusters appear to remain intact. This indicates that our models retain discriminative properties in the source domain, a fact we verify experimentally.

Maximum Mean Discrepancy (MMD): is based on the idea that two distributions are identical if and only if all their moments are identical. A divergence can be defined if we can measure how “different” the moments of the two distributions are. MMD is a method of efficiently doing this via

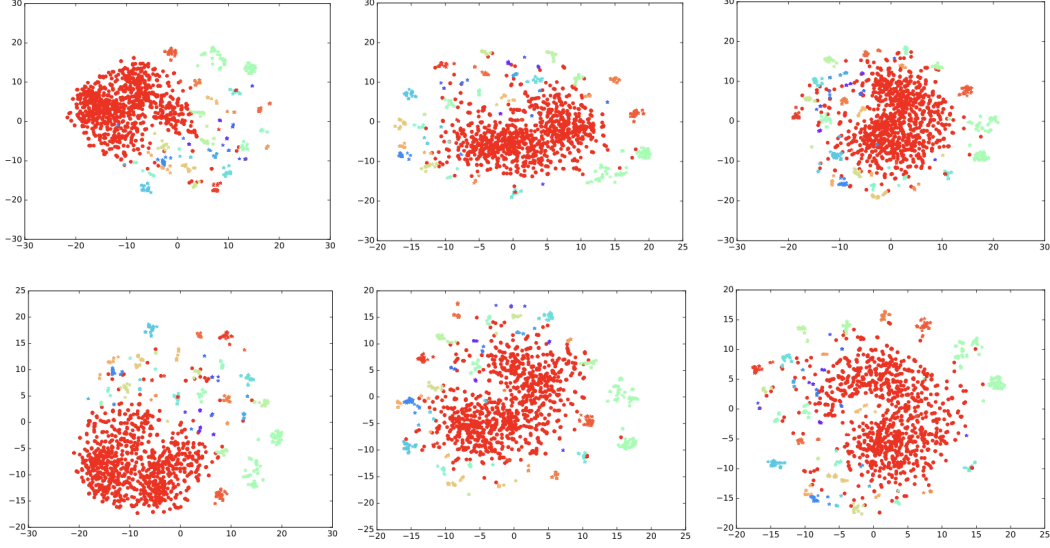


Figure 2: Domain Adversarial Neural Speaker Embedding Model

the kernel trick:

$$\begin{aligned} MMD(p(z)||q(z)) = \\ \mathbb{E}_{p(z),p(z')}[k(z,z')] + \mathbb{E}_{q(z),q(z')}[k(z,z')] - 2\mathbb{E}_{p(z),q(z')}[k(z,z')] \end{aligned} \quad (10)$$

In order to quantitatively evaluate our models in terms of domain adaptation, we measure the Maximum Mean Discrepancy distance between a selection of source data and the unlabelled target data. MMD is a standard distribution distance metric and has been applied in the context of domain adaptation [].

Fréchet Distance: The Fréchet Inception Distance (fid) is a popular approach for evaluating GANs, and has been shown to correlate well with human judgement of visual quality. Instead of an Inception network, we extract embeddings from our gan models from the source and target data. The Fréchet Distance between the Gaussian (m_s, C_s) obtained from the source data distribution p_s and the Gaussian (m_t, C_t) from the target data is given by:

$$d^2((\mathbf{m}_s, \mathbf{C}_s), (\mathbf{m}_t, \mathbf{C}_t)) = \|\mathbf{m}_s - \mathbf{m}_t\|_2^2 + Tr(\mathbf{C}_s + \mathbf{C}_t - 2(\mathbf{C}_s \mathbf{C}_t)^{1/2}) \quad (11)$$

Source Domain Speaker Verification: We use the same source data used to compute the MMD and Fréchet Distance to construct a trial list for verification. The list consists of 2500 recordings and we score them all versus all. There are a total of

From Fig. 3 we see that MMD and the Fréchet distance display similar trends. Surprisingly we see that Gradient Reversal only has a small effect on either metric, while the GAN models all have much lower MMD and Fréchet distances. We note that the model using the novel generator objective shows the lowest scores on both metrics. The results on source domain speaker verification also indicate that our models remain discriminative in the source domain as well, with only a small degradation as compared to the unadapted model. Interestingly, the Gradient Reversal model shows the best performance on this experiment albeit by a small margin.

7 Conclusion

In this work we presented a novel framework for learning domain-invariant speaker embeddings using GANs. By combining a powerful deep feature extractor, an end-to-end loss function and most importantly, adversarial training we are able to learn extremely compact speaker embeddings that deliver robust verification performance on challenging evaluation data. We showed that the

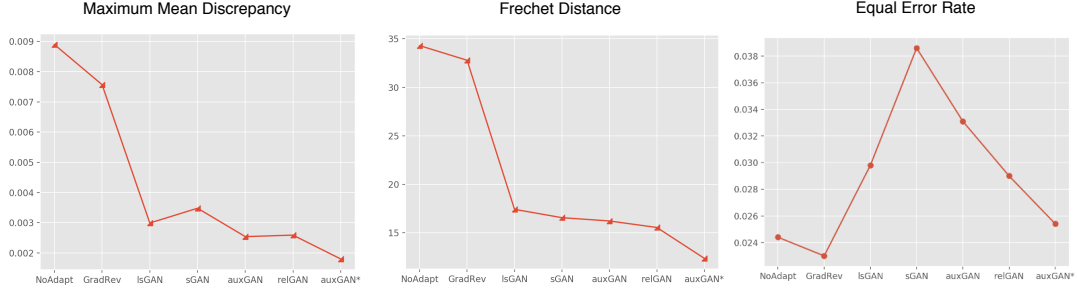


Figure 3: Domain Adversarial Neural Speaker Embedding Model

proposed methods do reduce the domain mismatch between source and target data in terms of MMD and Fréchet distance. Furthermore, we see that our methods adapt while maintaining their speaker discriminative nature in the source domain as well. In future work we will experiment with other GAN variants in an attempt to further improve performance.

References

- [1] Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations for domain adaptation. In *Advances in neural information processing systems*, pages 137–144, 2007.
- [2] Gautam Bhattacharya, Md Jahangir Alam, Vishwa Gupta, and Patrick Kenny. Deeply fused speaker embeddings for text-independent speaker verification. *Proc. Interspeech 2018*, pages 3588–3592, 2018.
- [3] Najim Dehak, Patrick J Kenny, Réda Dehak, Pierre Dumouchel, and Pierre Ouellet. Front-end factor analysis for speaker verification. *IEEE Transactions on Audio, Speech, and Language Processing*, 19(4):788–798, 2011.
- [4] Chao Li, Xiaokong Ma, Bing Jiang, Xiangang Li, Xuwei Zhang, Xiao Liu, Ying Cao, Ajay Kannan, and Zhenyao Zhu. Deep speaker: an end-to-end neural speaker embedding system. *arXiv preprint arXiv:1705.02304*, 2017.
- [5] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphreface: Deep hypersphere embedding for face recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 1, 2017.
- [6] Mitchell McLaren, Luciana Ferrer, Diego Castan, and Aaron Lawson. The speakers in the wild (sitw) speaker recognition database. In *Interspeech*, pages 818–822, 2016.
- [7] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier gans. *arXiv preprint arXiv:1610.09585*, 2016.
- [8] Seyed Omid Sadjadi, Timothée Kheyrkhan, Audrey Tong, Craig Greenberg, Elliot Singer Reynolds, Lisa Mason, and Jaime Hernandez-Cordero. The 2016 nist speaker recognition evaluation. In *Proc. Interspeech*, pages 1353–1357, 2017.
- [9] Swami Sankaranarayanan, Yogesh Balaji, Carlos D Castillo, and Rama Chellappa. Generate to adapt: Aligning domains using generative adversarial networks. *ArXiv e-prints, abs/1704.01705*, 2017.
- [10] Jian Shen, Yanru Qu, Weinan Zhang, and Yong Yu. Wasserstein distance guided representation learning for domain adaptation. In *AAAI*, 2018.
- [11] David Snyder, Daniel Garcia-Romero, Gregory Sell, Daniel Povey, and Sanjeev Khudanpur. X-vectors: Robust dnn embeddings for speaker recognition. *ICASSP, Calgary*, 2018.
- [12] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 4, 2017.

- 272 [13] Li Wan, Quan Wang, Alan Papir, and Ignacio Lopez Moreno. Generalized end-to-end loss for
273 speaker verification. In *2018 IEEE International Conference on Acoustics, Speech and Signal*
274 *Processing (ICASSP)*, pages 4879–4883. IEEE, 2018.
- 275 [14] Feng Wang, Jian Cheng, Weiyang Liu, and Haijun Liu. Additive margin softmax for face
276 verification. *IEEE Signal Processing Letters*, 25(7):926–930, 2018.
- 277 [15] Yingke Zhu, Tom Ko, David Snyder, Brian Mak, and Daniel Povey. Self-attentive speaker
278 embeddings for text-independent speaker verification. *Proc. Interspeech 2018*, pages 3573–3577,
279 2018.