

Xi-Vector Embedding for Speaker Recognition

Kong Aik Lee, *Senior Member, IEEE*, Qionqiong Wang, and Takafumi Koshinaka

Abstract—We present a Bayesian formulation for deep speaker embedding, wherein the xi-vector is the Bayesian counterpart of the x-vector, taking into account the uncertainty estimate. On the technology front, we offer a simple and straightforward extension to the now widely used x-vector. It consists of an auxiliary neural net predicting the frame-wise uncertainty of the input sequence. We show that the proposed extension leads to substantial improvement across all operating points, with a significant reduction in error rates and detection cost. On the theoretical front, our proposal integrates the Bayesian formulation of linear Gaussian model to speaker-embedding neural networks via the pooling layer. In one sense, our proposal integrates the Bayesian formulation of the i-vector to that of the x-vector. Hence, we refer to the embedding as the xi-vector, which is pronounced as /zai/ vector. Experimental results on the SITW evaluation set show a consistent improvement of over 17.5% in equal-error-rate and 10.9% in minimum detection cost.

Index Terms—Speaker Verification, Neural Embedding, Uncertainty

I. INTRODUCTION

Automatic speaker recognition is the task of identifying or verifying an individual’s identity from samples of his/her voice using machine learning algorithms, without any human intervention [1]. How speaker recognition is carried out has changed substantially throughout these years with increasing accuracy and robustness against various sources of variability in the speech signal. Modern speaker recognition systems [2], [3], [4], [5], [6] consist of a speaker embedding front-end followed by a scoring backend. In this approach, speech utterances are first represented as fixed-length vectors – the so-called speaker embeddings. The current de-facto standard of speaker embedding is x-vector [7]. For scoring backend, probabilistic linear discrimination analysis (PLDA) [8], [9] is commonly used.

Unlike passwords that have zero uncertainty conditioned on a person’s identity [10], human voices exhibit both extrinsic and intrinsic variability. Major sources of extrinsic variability in speech signals are background noises, channel distortion, and room acoustics. Intrinsic factors include the physiological nature of the vocal apparatus and psychological states (e.g., emotion, mental health condition, etc.) of the speaker, and the biological constraint of the vocal tract leads to acoustically different utterances every time we repeat the same sentence. X-vectors, and similar forms of deep speaker embedding, do not consider the uncertainty of features. In a restricted

sense, uncertainty is merely captured implicitly with empirical variance estimates at the utterance level. Consequently, they show low robustness against local and random perturbation which is the inherent property of speech utterances.

The ability to handle uncertainty has been the cornerstone in the successful use of generative models for speaker recognition [11], [12], [13]. In the *universal background model* or UBM, the uncertainty of feature vectors is modeled with the covariance matrices associated with its Gaussian components [11]. In the i-vector embedding, uncertainty is captured with the prior and posterior covariance matrices of its latent variable [12], [13].

Under the big data regime, deep learning is used to achieve state-of-the-art performance, notwithstanding that most speaker-embedding neural networks are not able to represent uncertainty. In this paper, we bridge the gap by incorporating uncertainty modeling in a speaker embedding neural network. To this end, we first estimate the frame-wise uncertainty of the input sequences with an auxiliary neural network. This operation is handled with a linear Gaussian model at the pooling layer and trained as part of the neural network.

In addition to good performance, our study sheds light on the role of uncertainty in speaker embedding. It turns out that the variance estimate is used as the indicator to which feature vectors, and which dimensions of the feature vectors, are useful in speaker classification. Higher weights are given to those frames and features with lower uncertainty and vice versa. In a way, it plays a role similar to the attention model reported in [14], [15]. The difference here is that the usefulness, or importance of feature vectors, is associated with the uncertainty estimate while this is loosely defined in prior works. Furthermore, the prior in the linear Gaussian model has its effect on the weights assigned to each frame and feature. More importantly, we show that a generative model could be inserted as part of a discriminative neural network to handle uncertainty.

II. NEURAL SPEAKER EMBEDDINGS

We define speaker embedding vectors as the representation of variable-length utterances as fixed-length continuous-valued vectors. Embeddings from the same speakers are close together in the embedding vector space, and therefore allow easy comparison between speakers with simple geometric operations. Depending on how the embedding extractor is trained, we divide speaker embeddings into two broad categories, namely, (i) unsupervised embedding, and (ii) supervised embedding. The celebrated i-vector [13], and the classical *Gaussian mixture model* (GMM) supervector [16] and alike [17], belong to the first category. A popular example of the second category is the x-vector embeddings [7], [18]. Different from that of

K. A. Lee was with the Biometrics Research Laboratories, NEC Corporation, Japan. He is now with the Institute for Infocomm Research, A*STAR, Singapore (e-mail: lee_kong_aik@i2r.a-star.edu.sg).

Q. Wang is with the Biometrics Research Laboratories, NEC Corporation, Japan (e-mail: q-wang@nec.com).

K. Takafumi was with the Biometrics Research Laboratories, NEC Corporation, Japan. He is now with the School of Data Science, Yokohama City University, Japan (e-mail: koshinak@yokohama-cu.ac.jp).

the latter, both i-vector and GMM supervector are based on a generative model trained using an unsupervised maximum likelihood criterion. Supervised embeddings rely on the use of labelled data for discriminative training, typically, with a multi-class speaker-discriminative loss.

Recent development has shown the benefit of supervised learning with deep neural networks (DNNs) coupled with the massive use of data augmentation. Both have advanced the performance of neural speaker embedding by a large margin. The conventional x-vector extractor is a DNN consisting of three functional blocks:

- An **encoder** (or a frame processor) implemented with multiple layers of time-delay neural network (TDNN) [19], [20]. In [21], it's shown that a TDNN could be implemented as a 1D-CNN with dilation. In [4], [22], [23], 2D-convolution has shown to be effective as well.
- A **temporal pooling layer** to compute an aggregated measure from the frame-level feature vectors produced by the encoder.
- A **decoder** to classify the utterances to speaker classes. One of the layers is designed to be a bottleneck layer, the output of which (after affine projection, and before the non-linearity) is the so-called x-vector speaker embedding.

One element that is missing in the current framework is the ability to handle uncertainty induced by the random factors inherent in the generation (intrinsic) and transmission (extrinsic) of human voices. In particular, the frame-level features produced by the encoder are point estimates, which do not consider the uncertainty of the latent representation. In this paper, we bridge the gap with the provision of (i.) an encoder that infers the uncertainty in addition to the point estimate of frame-level features, and (ii.) a generative model that leverages the uncertainty measure in deep speaker embedding. It is worth mentioning that the frame-wise uncertainty estimated by the encoder neural network is regarded as *aleatoric* uncertainty [24] in Bayesian deep learning.

III. UNCERTAINTY MODELING IN THE LATENT SPACE

To incorporate uncertainty measurement and modeling into the neural speaker embedding framework, we introduce two new concepts here.

A. Uncertainty estimation

The encoder maps an input frame \mathbf{x}_t to a point estimate \mathbf{z}_t in a latent space. We propose to characterize the frame uncertainty with a covariance matrix \mathbf{L}_t^{-1} associated with each estimate \mathbf{z}_t . This is accomplished by using a base neural network with two heads, one for \mathbf{z}_t and one for \mathbf{L}_t^{-1} , and train simultaneously. Without loss of generality, we denote these operations as

$$\mathbf{z}_t = f_{\text{enc}}(\mathbf{x}_t | \mathbf{x}_t^{t \pm n}) \quad \text{and} \quad \log \mathbf{L}_t = g_{\text{enc}}(\mathbf{x}_t | \mathbf{x}_t^{t \pm n}) \quad (1)$$

where a context of $\pm n$ neighbouring frames is taken into account for each estimate. Note that the second equation in (1) is absent in the conventional x-vector embedding. Also,

we assume that the covariance (and precision) matrices are diagonal and chose to estimate directly the log-precision which turns out to be more convenient as we shall see in the next section.

With the construct in (1), the frame encoder maps an input sequence $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$ of length T to another sequence of enhanced features $\{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_T\}$ with their corresponding uncertainty measures given by $-1 \times \log \{\mathbf{L}_1, \mathbf{L}_2, \dots, \mathbf{L}_T\}$.

B. Posterior inference using frame-wise uncertainty

We assume that a *linear Gaussian model* is responsible for generating the representations \mathbf{z}_t , as follows:

$$\begin{aligned} \text{Generative model:} \quad & \mathbf{z}_t = \mathbf{h} + \epsilon_t \\ \text{Latent variable:} \quad & \mathbf{h} \sim \mathcal{N}(\mu_p, \mathbf{L}_p^{-1}) \\ \text{Uncertainty:} \quad & \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{L}_t^{-1}) \end{aligned} \quad (2)$$

Here, \mathbf{h} is the latent variable assigned to the entire sequence, and ϵ_t is a random variable caters for the uncertainty covariance measure for each frame. We also impose a Gaussian prior on the variable \mathbf{h} with the prior mean vector μ_p and covariance matrix \mathbf{L}_p^{-1} .

Given the input sequence and uncertainty estimate, it can be shown that the posterior distribution of \mathbf{h} is also Gaussian:

$$p(\mathbf{h} | \mathbf{z}_1, \dots, \mathbf{z}_T, \mathbf{L}_1^{-1}, \dots, \mathbf{L}_T^{-1}) = \mathcal{N}(\mathbf{h} | \phi_s, \mathbf{L}_s^{-1}) \quad (3)$$

with its posterior mean vector and precision matrix computed as

$$\phi_s = \mathbf{L}_s^{-1} \left[\sum_{t=1}^T \mathbf{L}_t \mathbf{z}_t + \mathbf{L}_p \mu_p \right] = \sum_{t=0}^T \mathbf{A}_t \mathbf{z}_t \quad (4)$$

and

$$\mathbf{L}_s = \sum_{t=1}^T \mathbf{L}_t + \mathbf{L}_p = \sum_{t=0}^T \mathbf{L}_t \quad (5)$$

respectively. For sake of clarity, we have assigned the index $t = 0$ to the prior such that $\mathbf{z}_0 = \mu_p$, $\mathbf{L}_0 = \mathbf{L}_p$, and therefore

$$\mathbf{A}_t = \mathbf{L}_s^{-1} \mathbf{L}_t \quad \text{for } t = 0, 1, \dots, T \quad (6)$$

Notice that the posterior mean in (4) is the weighted average of latent feature vectors \mathbf{z}_t and prior μ_p , each weighted with a gain factor \mathbf{A}_t determined by the uncertainty measure \mathbf{L}_t^{-1} and the posterior covariance \mathbf{L}_s^{-1} .

The posterior inference entails a temporal aggregation operation since we use one latent variable \mathbf{h} for the entire sequence. A similar operation was used in the classical i-vector inference. The difference lies in how the uncertainty measures are derived. In the i-vector paradigm, the uncertainty covariance matrix is drawn from a finite set (i.e., the covariance matrices of the UBM). Here, the uncertainty is predicted on the fly using a neural network. We refer to this as the *heteroscedastic aleatoric* uncertainty.

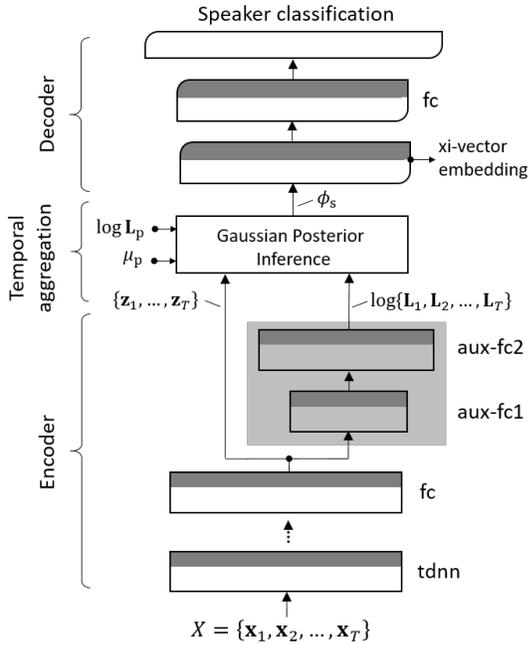


Fig. 1. An exemplary xi-vector embedding neural network implemented with multiple TDNN and fully-connected (fc) layers. The round corner indicates layers that process utterance-level representations at the decoder. Shaded bars indicate non-linear activation functions. All layers use ReLU except for aux-fc2 that uses a *softplus* followed by $2 \times \log$ operation for log-precision estimation and a *softmax* at the output.

IV. XI-VECTOR EMBEDDING WITH UNCERTAINTY

We aim to incorporate frame uncertainty measures into deep speaker representation learning. This is achieved by inserting the generative linear Gaussian model into the speaker embedding neural network. Our proposal is based on the observation that Gaussian posterior inference is essentially a mapping operation where sequences of point estimates, \mathbf{z}_t , and uncertainty measures, \mathbf{L}_t^{-1} , are mapped to a fixed-length posterior mean vector. Figure 1 shows an exemplary neural network that implements the model. Compared to the conventional x-vector, three new components are:

- Frame uncertainty is inferred at the output of the frame encoder in addition to a point estimate,
- Gaussian posterior inference is used for temporal aggregation, and
- Posterior mean vector is used as input to the decoder replacing the first and second-order moments in conventional x-vector.

The posterior inference entails a temporal aggregation operation replacing the simplistic statistical pooling used in the conventional x-vector. The second-order moment is no longer required for utterance classification layers since frame uncertainty has been accounted for. As in the x-vector framework, an embedding representation is obtained by taking the pre-activation output of the first hidden layer of the decoder network.

Algorithm 1 shows the minibatch stochastic gradient descent training. In lines 6 and 7, the encoder produces a point estimate and the log-precision for each input frame, respectively. In

Algorithm 1: Stochastic gradient descent training of the xi-vector neural network. A training example consists of a speech segment X_s of T frames and a speaker label y_s . The set θ includes network parameters, the prior mean, and log-precision.

```

1 Initialize:  $\mu_p \leftarrow \mathbf{0}$ ,  $\log \mathbf{L}_p \leftarrow \mathbf{0}$ ;
2 while stopping criterion not met do
3   Sample a mini-batch of  $M$  examples  $(X_s, y_s)_{s=1}^M$ ;
4    $\mathbf{g} \leftarrow \mathbf{0}$ ;
5   for  $s = 1$  to  $M$  do
6      $\{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_T\} \leftarrow f_{\text{enc}}(X_s)$ ;
7      $\log \{\mathbf{L}_1, \mathbf{L}_2, \dots, \mathbf{L}_T\} \leftarrow g_{\text{enc}}(X_s)$ ;
8      $\log \mathbf{L}[i] \leftarrow \log \{\mathbf{L}_0[i], \mathbf{L}_1[i], \dots, \mathbf{L}_T[i]\}$ ;
9      $\mathbf{A}[i] \leftarrow \text{softmax}(\log \mathbf{L}[i])$ ,  $i = 1, 2, \dots, D$ ;
10     $\phi_s \leftarrow \sum_{t=0}^T \mathbf{A}_t \mathbf{z}_t$ ;
11     $\hat{y}_s \leftarrow f_{\text{dec}}(\phi_s)$ ;
12    Compute loss  $L(\hat{y}_s, y_s)$ ;
13     $\mathbf{g} \leftarrow \mathbf{g} + \frac{1}{M} \nabla L$ ;
14  end
15   $\theta \leftarrow \theta + \eta \mathbf{g}$ 
end

```

lines 8 and 9, the algorithm estimates the gain factor using (5) and (6). Since we assume that the precision matrix is diagonal and the encoder predicts directly the log-precision, the i -th diagonal element of the gain factor is computed as:

$$\mathbf{A}_t[i] = \frac{\mathbf{L}_t[i]}{\sum_{t'=0}^T \mathbf{L}_{t'}[i]} = \frac{e^{\log \mathbf{L}_t[i]}}{\sum_{t'=0}^T e^{\log \mathbf{L}_{t'}[i]}} \quad (7)$$

for $i = 1, \dots, D$, where D is the dimension of the pooling layer. Line 9 implements Eq. (7) with a *softmax* function across the temporal index t . The entire neural network is trained by minimizing the cross-entropy loss.

Note that the prior mean and log-precision are initialized to zeros. The prior parameters are updated simultaneously with other network parameters via a posterior inference step followed by gradient back-propagation, which combines the generative nature of the i-vector with the benefit of speaker discrimination of the x-vector. Hence, we refer to the speaker embeddings as the xi-vectors – pronounced as the /zai/ vector.

V. EXPERIMENTS

We present two sets of experiments. The first set was conducted on the Speakers in the Wild (SITW) *dev* and *eval* sets [25]. The training data was drawn from the VoxCeleb-1 and VoxCeleb-2 corpora [5]. The second set of experiments was conducted on the NIST SRE'18 and SRE'19 *eval* sets. The training set consists primarily of English speech corpora, which encompasses Switchboard, Fisher, and the MIXER corpora used in SREs 04 – 06, 08, and 10. Since the SRE eval sets consist of enrollment and test segments in Tunisian Arabic, domain adaptation was performed on the PLDA using the unlabeled subsets provided for the evaluation. We used 40 and 23-dimensional MFCCs with 10ms frameshift for the first and second sets of experiments, respectively. Mean-normalization

TABLE I
PERFORMANCE COMPARISON OF XI-VECTOR EMBEDDING WITH
X-VECTOR ON THE SITW DEV AND EVAL SETS.

	SITW-Dev		SITW-Eval	
	EER (%)	MinDCF	EER (%)	MinDCF
x-vector (μ, σ)	2.38	0.251	2.65	0.278
xi-vector (ϕ, σ)	1.72	0.216	2.30	0.239
x-vector (μ)	2.93	0.333	3.23	0.362
xi-vector (ϕ)	1.96	0.214	2.19	0.248

TABLE II
ABLATION ANALYSIS OF XI-VECTOR ON SITW DEV AND EVAL SETS.

	SITW-Dev		SITW-Eval	
	EER (%)	MinDCF	EER (%)	MinDCF
No-Prior (ϕ)	1.96	0.235	2.43	0.268
Isotropic (ϕ)	2.52	0.273	3.03	0.325
No-Prior (ϕ, σ)	2.10	0.240	2.35	0.254
Isotropic (ϕ, σ)	2.28	0.244	2.54	0.264

over a sliding window of 3s and energy-based VAD were then applied. Data augmentation [26] was performed on the training sets using the MUSAN dataset [27]. As in most state-of-the-art implementations, speaker embeddings were reduced to 200 dimensions via LDA projection before PLDA.

We first compare the performance of the proposed xi-vector to the conventional x-vector baseline. We used a TDNN-5 base architecture [7], [14] for both x-vector and xi-vector extraction, with a pooling layer of $D = 1500$ dimensions for all experiments. For the xi-vectors, frame precision matrices were assumed to be diagonal and estimated with a two-layer feed-forward neural network ($1500 \times 256 \rightarrow 256 \times 1500$). The size of the hidden layer was set to 256 so that the number of parameters of the auxiliary neural network amounts to $(1500 \times 256) \times 2$, which is the same as the number of weights required to map the standard deviation of the pooled statistics to form the x-vector, i.e., 1500×512 . Note that the standard deviation is no longer required in the xi-vector.

Table I shows the performance comparison in terms of the EER and MinDCF ($P_{target} = 0.01$). Comparing the baseline x-vector (μ, σ) to the xi-vector (ϕ) in the first and last rows of Table I, respectively, the results show that the proposed xi-vector embedding gives consistent improvement on both EER and MinDCF. The performance gain amounts to 17.5% in EER and 10.9% in MinDCF on the SITW-Eval set. Similar performance gain could be observed on SITW-Dev set. The results in the second and third rows of Table I show that the use of standard deviation is not required in the xi-vector but essential for x-vector embeddings. This proves the case that frame uncertainty estimate gives a better account for variability in the inputs. Note that we used the gain factor in (7) to compute the weighted standard deviation for the xi-vector results in the second last row of Table I.

Table II show the results of ablation analysis. We removed the prior (denoted as *No-Prior*) and further replaced the diagonal variance estimate with a spherical covariance estimate (denoted as *No-Prior + Isotropic*). Comparing the results in the first two rows to that of the xi-vector (ϕ) in Table I, we observe consistent degradation across all decision points

TABLE III
PERFORMANCE COMPARISON OF XI-VECTOR EXTRACTOR WITH
X-VECTOR.

	SRE'18-Eval		SRE'19-Eval	
	EER (%)	MinDCF	EER (%)	MinDCF
x-vector (μ, σ)	7.53	0.503	6.97	0.510
xi-vector (ϕ)	7.13	0.483	6.34	0.478

except the EER on the dev set. On the eval set, the EER degradation amounts to 11.2% with the prior removed and further to 38.7% when the isotropic variance assumption was imposed. The results in the last two rows of Table II corresponds to the xi-vector (ϕ, σ) in Table I. Similar degradation could be observed with the prior removed and the isotropic assumption, though the degradation is considerably smaller. Comparing the last two row of Table II to the results with x-vector (ϕ, σ) in Table I, it is clear that inclusion of frame uncertainty estimate helps. Above all, these results show that a proper frame uncertainty estimation and the use of Gaussian posterior inference formulation for temporal aggregation outperform simplistic statistical pooling which depends on the empirical mean and standard deviation estimate.

The second set of experiments was conducted on the SRE'18 and SRE'19 eval sets. Table III shows the performance comparison in terms of EER and MinDCF ($P_{target} = 0.01$). We observe consistent performance improvement with xi-vector on both test sets. The performance gain on the SRE'18 is 5.3% in EER and 4.0% in MinDCF. The performance gain is slightly higher on SRE'19 which amounts to 9.0% and 6.3% in terms of EER and MinDCF, respectively. These results confirm the effectiveness of the proposed xi-vector on the narrowband SRE and wideband SITW test sets.

VI. CONCLUSION

From i-vector to x-vector, and xi-vector, the central theme is speaker representation learning – to find fixed-length representations that allow easy comparison between speakers with simple geometric operations. What sets apart the proposed xi-vector from its predecessor is the use of generative modeling in a supervised training framework. One essential element of a generative model is its ability to handle uncertainty. Discriminative embeddings, like x-vectors, do not take the uncertainty (distribution) of features into consideration. We propose to characterize the uncertainty of input sequences with frame-wise uncertainty estimates, which are then used with the point estimates to derive speaker embedding vectors. This is accomplished with a linear Gaussian model trained as part of the embedding neural network. Given the capability to take into account frame-wise uncertainty, the proposed xi-vector embeddings exhibit improve robustness to perturbation and therefore better performance.

VII. ACKNOWLEDGEMENTS

The authors would like to thank Mr. Hitoshi Yamamoto and Dr. Shinya Miyakawa at NEC Biometrics Research Laboratories for giving the opportunity to begin and complete this work.

REFERENCES

- [1] K. A. Lee, O. Sadjadi, H. Li, and D. Reynolds, "Two decades into speaker recognition evaluation - are we there yet?," *Computer Speech & Language*, vol. 61, p. 101058, 2020.
- [2] K. A. Lee, V. Hautamaki, T. Kinnunen, H. Yamamoto, K. Okabe, V. Vestman, J. Huang, G. Ding, H. Sun, A. Larcher, R. K. Das, H. Li, M. Rouvier, P.-M. Bousquet, W. Rao, Q. Wang, C. Zhang, F. Bahmaninezhad, H. Delgado, and M. Todisco, "I4U submission to NIST SRE 2018: Leveraging from a decade of shared experiences," in *Proc. Interspeech*, pp. 1497–1501, 2019.
- [3] P. Matejka, O. Plchot, O. Glembek, L. Burget, J. Rohdin, H. Zeinali, L. Mosner, A. Silnova, O. Novotny, M. Diez, and J. Černocky, "13 years of speaker recognition research at BUT, with longitudinal analysis of nist sre," *Computer Speech & Language*, vol. 63, p. 101035, 2020.
- [4] J. Villalba, N. Chen, D. Snyder, D. Garcia-Romero, A. McCree, G. Sell, J. Borgstrom, L. P. Garcia-Perera, F. Richardson, R. Dehak, P. A. Torres-Carrasquillo, and N. Dehak, "State-of-the-art speaker recognition with neural network embeddings in NIST SRE18 and Speakers in the Wild evaluations," *Computer Speech & Language*, vol. 60, p. 101026, 2020.
- [5] A. Nagrani, J. S. Chung, W. Xie, and A. Zisserman, "Voxceleb: Large-scale speaker verification in the wild," *Computer Speech & Language*, vol. 60, p. 101027, 2020.
- [6] K. A. Lee, H. Yamamoto, K. Okabe, Q. Wang, L. Guo, T. Koshinaka, J. Zhang, and K. Shinoda, "NEC-TT system for mixed-bandwidth and multi-domain speaker recognition," *Computer Speech & Language*, vol. 61, p. 101033, 2020.
- [7] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, "X-vectors: Robust DNN embeddings for speaker recognition," in *Proc. ICASSP*, pp. 5329–5333, 2018.
- [8] S. Ioffe, "Probabilistic linear discriminant analysis," in *proceedings of the 9th European Conference on Computer Vision*, 2006.
- [9] S. J. D. Prince and J. H. Elder, "Probabilistic linear discriminant analysis for inferences about identity," in *Proc. ICCV*, pp. 1–8, 2007.
- [10] K. Takahashi and T. Murakami, "A measure of information gained through biometric systems," *Image and Vision Computing*, vol. 32, no. 12, pp. 1194–1203, 2014.
- [11] D. A. Reynolds, T. F. Quatieri, and R. B. Dunn, "Speaker verification using adapted gaussian mixture models," *Digital Signal Processing*, pp. 19–41, 2000.
- [12] P. Kenny, P. Ouellet, N. Dehak, V. Gupta, and P. Dumouchel, "A study of interspeaker variability in speaker verification," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 16, no. 5, pp. 980–988, 2008.
- [13] N. Dehak, P. 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.
- [14] K. Okabe, T. Koshinaka, and K. Shinoda, "Attentive statistics pooling for deep speaker embedding," in *Proc. Interspeech*, pp. 2252–2256, 2018.
- [15] Q. Wang, K. Okabe, K. A. Lee, H. Yamamoto, and T. Koshinaka, "Attention mechanism in speaker recognition: What does it learn in deep speaker embedding?," in *Proc. IEEE SLT Workshop*, pp. 1052–1059, 2018.
- [16] W. M. Campbell, D. E. Sturim, and D. A. Reynolds, "Support vector machines using gmm supervectors for speaker verification," *IEEE Signal Processing Letters*, vol. 13, no. 5, pp. 308–311, 2006.
- [17] K. A. Lee, C. You, H. Li, T. Kinnunen, and K. C. Sim, "Using discrete probabilities with bhattacharyya measure for svm-based speaker verification," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 4, pp. 61–870, 2011.
- [18] B. Desplanques, J. Thienpondt, and K. Demuynck, "ECAPA-TDNN: Emphasized channel attention, propagation and aggregation in TDNN based speaker verification," in *Proc. Interspeech 2020*, pp. 3830–3834, 2020.
- [19] A. Waibel, T. Hanazawa, G. Hinton, K. Shikano, and K. J. Lang, "Phoneme recognition using time-delay neural networks," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 37, no. 3, pp. 328–339, 1989.
- [20] V. Peddinti, D. Povey, and S. Khudanpur, "A time delay neural network architecture for efficient modeling of long temporal contexts," in *Interspeech 2015*, pp. 3214–3218, 2015.
- [21] S. Bai, J. Z. Kolter, and V. Koltun, "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling," *CoRR*, vol. abs/1803.01271, 2018.
- [22] H. Zeinali, S. Wang, A. Silnova, P. Matějka, and O. Plchot, "BUT system description to voxceleb speaker recognition challenge 2019," 2019.
- [23] J. S. Chung, A. Nagrani, and A. Zisserman, "Voxceleb2: Deep speaker recognition," in *Proc. Interspeech 2018*, pp. 1086–1090, 2018.
- [24] A. Kendall and Y. Gal, "What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision?," in *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS'17*, p. 5580–5590, 2017.
- [25] M. McLaren, L. Ferrer, D. Castan, and A. Lawson, "The Speakers in the Wild (SITW) speaker recognition database," in *Proc. Interspeech*, pp. 818–822, 2016.
- [26] T. Ko, V. Peddinti, M. S. Daniel Povey, and S. Khudanpur, "A study on data augmentation of reverberant speech for robust speech recognition," in *Proc. IEEE ICASSP*, pp. 5220–5224, 2017.
- [27] D. Snyder, G. Chen, and D. Povey, "MUSAN: a music, speech, and noise corpus," in *arXiv:1510.08484*, 2015.