

DOMAIN MISMATCH ROBUST ACOUSTIC SCENE CLASSIFICATION USING CHANNEL INFORMATION CONVERSION

Seongkyu Mun¹, Suwon Shon²

Clova AI Research, Naver Corp. Seongnam, South Korea¹
MIT Computer Science and Artificial Intelligence Laboratory, Cambridge, MA, USA²

sk.moon@navercorp.com

swshon@csail.mit.edu

ABSTRACT

In a recent acoustic scene classification (ASC) research field, training and test device channel mismatch have become an issue for the real world implementation. To address the issue, this paper proposes a channel domain conversion using factorized hierarchical variational autoencoder. Proposed method adapts both the source and target domain to a pre-defined specific domain. Unlike the conventional approach, the relationship between the target and source domain and information of each domain are not required in the adaptation process. Based on the experimental results using the IEEE detection and classification of acoustic scenes and event 2018 task 1-B dataset and the baseline system, it is shown that the proposed approach can mitigate the channel mismatching issue of different recording devices.

Index Terms— acoustic scene classification, factorized hierarchical variational autoencoder, domain adaptation

1. INTRODUCTION

Acoustic Scene Classification (ASC) is a task that classifies input sounds into specific acoustic scenes, such as office, park, airport, tram, etc. In previous ASC researches, transfer learning [1, 2], attention mechanism [3, 4] and DB augmentation [5, 6] were proposed to improve ASC performance. Recently, researches on the ASC have been intensively studied on the IEEE Detection and Classification of Acoustic Scenes and Events (DCASE) 2016-2018 challenges [7, 8, 9]. Due to the simple and clear task of classifying pre-defined scene labels for sound data of specific length, various techniques from acoustic signal processing fields, such as speaker recognition and music information retriever, have been tried. All of the top teams in the last three years used Convolutional Neural Network (CNN)-based architectures and additionally used I-vectors [10], Generative Adversarial Networks (GAN) based DB augmentation [5] and harmonic-percussive source separation based pre-processing [11, 12], respectively. Unlike 2016 and 2017, the 2018 DCASE challenge task 1 added sub-task B, which addresses the dataset recorded with different

devices [9]. In real-world environments, the device mismatching issue is inevitable, so the new subtask has a practically important issue. The task consists of a relatively large source domain A (recorded by device A) and a relatively small target domain B and C (recorded by device B and C, respectively). Submissions of the challenge task 1-B are ranked by classification accuracy of target devices B and C. To the best of our knowledge, during the challenge period, there were no submitted technical reports that handle different device issue directly. After the challenge, in order to address the related issue, a paper using GAN based domain adaptation has been submitted for the following DCASE 2018 workshop [13]. The adaptation module (from ‘target’ to ‘source’ domain) and domain discriminator are optimized through adversarial training, and the ASC performance was improved on the DCASE 2018 task 1-B DB.

Although the aforementioned approach effectively adapted the device domain, there is a limitation that DB of the target domain is required for training the adaption module. This limitation could be critical in some cases. For example, when input sounds are received via other unseen devices or web-streaming, it is difficult to gather sufficient target domain DB for training the domain adaptation module.

To address the issue, this paper proposes an adaptation from *target or source to other specific* domain, not target to source domain or vice versa. Proposed method adapts the source domain as well as the target domain. As shown in Figure 1, the device domain (channel) related component (z_2) is disentangled from the input signal, then it is shifted to the other domain (e.g. universal domain), and the reconstruction process is followed. Since z_2 components of the input features are mapped to the specific domain, the relationship between the target and source domain and information of each domain are not required in the adaptation process. In order to implement the aforementioned process, we utilized Factorized Hierarchical Variational AutoEncoder (FHVAE) [14], which shows notable performance improvement in voice conversion and sequential information representation [15, 16, 17]. We adapted input features by using FHVAE for generating factors of the acoustic scene and device-related component and

mapping the device related component to the other specific domain. Based on the experimental results using the DCASE 2018 task 1-B dataset and the baseline system, it is shown that the proposed approach can mitigate the channel mismatching issue of different recording devices.

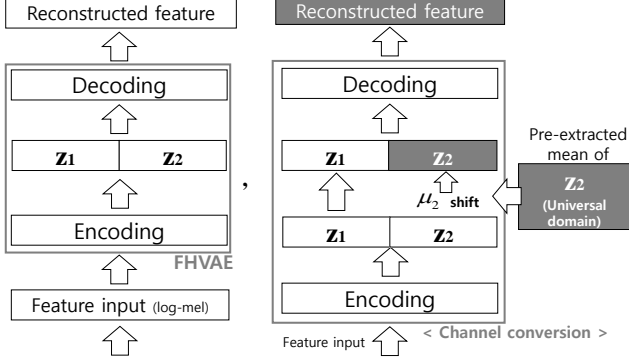


Fig. 1: The structure of FHVAE (Left) and concept of channel conversion (Right)

2. ACOUSTIC SCENE CLASSIFICATION IN CHANNEL MISMATCHED CONDITION

2.1. Channel related component disentanglement using FHVAE

In this section, we briefly describe the FHVAE for our ASC system. More details of FHVAE can be found in [14, 15, 16]. The FHVAE [14] is a variant of variational autoencoder [18] that models a probabilistic hierarchical generative process of sequential data, and learns disentangled and interpretable representations. Generation of a sequence of N segments involves one sequence-level latent variable, μ_2 , and N pairs of segment-level latent variable z_1 and z_2 . μ_2 , the prior component of sequence-dependent is drawn from $p(\mu_2) = N(\mu_2|0, \sigma_{\mu_2}^2 I)$. N i.i.d latent segment variables $Z_1 = \{z_1^{(n)}\}_{n=1}^N$ are drawn from a global prior $p(z_1) = N(z_1|0, \sigma_{z_1}^2 I)$. N i.i.d. latent sequence variables $Z_2 = \{z_2^{(n)}\}_{n=1}^N$ are drawn from a sequence-dependent prior $p(z_2|\mu_2) = N(z_2|\mu_2, \sigma_{z_2}^2 I)$. At last, N i.i.d. subsequences $X = \{x^{(n)}\}_{n=1}^N$ are drawn from $p(x|z_1, z_2) = \mathcal{N}(x|f_{\mu_x}(z_1, z_2), \text{diag}(f_{\sigma_x^2}(z_1, z_2)))$, where $f_{\mu_x}(\cdot, \cdot)$ and $f_{\sigma_x^2}(\cdot, \cdot)$ are parameterized by a decoder neural network. [17, 19]. Since the exact posterior inference is intractable, FHVAEs introduce an inference model to approximate the true posterior. In this work, we followed the latest training method of the FHVAE research [17]. Figure 2 shows the aforementioned model.

By imposing a sequence-dependent prior to z_2 , the model is encouraged to represent with z_2 the generating factors that are relatively consistent within a sequence. For example, such factors can include microphone frequency response and room

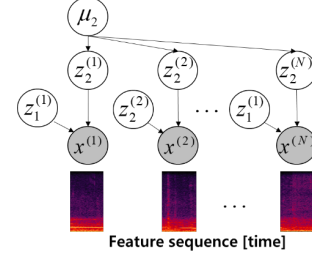


Fig. 2: Graphical illustration of the FHVAE generative model. Grey nodes denotes observed variables, and white nodes are the latent variables

impulse response. On the other hand, z_1 tends to encode information about the residual generating factors that change from segment to segment, such as acoustic scene related audio events. In order to compare the characteristics of the two latent variables, 2-dimension t -Distributed Stochastic Neighbor Embedding (t-SNE) projection examples of z_1 and z_2 distributions are shown in Figure 3. Each point represents one segment. The FHVAE model was trained using DCASE 2018 task 1-B development train set DB. Note that the device label information was not used in the FHVAE training process. Detail configurations of the model and inference will be discussed in section 3. As shown in Figure 3, z_1 shows channel invariant characteristic compared to z_2 , and z_2 distribution has a tendency to be clustered by each channel subset. The result of z_2 plot shows that blue dots (Tram, Device A) are closer to orange dots (Bus, Device A) of the same recording device, compared to the green dot (Tram, Device B) of the same acoustic scene class.

Assuming that the latent variable z_2 contains channel discriminative information, we propose a process of shifting z_2 values for equalizing (adapting) channel components. In addition, based on previous researches of channel invariant feature representation [19, 20], we conducted an experiment of using z_2 as a feature input of the classifier (without reconstruction) for performance comparison.

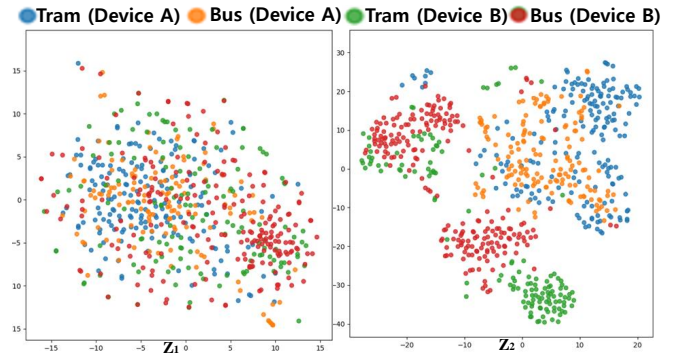


Fig. 3: Scatter plots of t-SNE projected z_1 and z_2 with models trained on DCASE 2018 task 1-B

2.2. Channel conversion using latent variable shifting

To utilize latent variable z_2 to convert channel, we obtained motivation from the previous research of [14]. For transforming sequence-level attributes while preserving segment-level attributes, we conducted the mean μ_2 shifting process. Based on the FHVAE framework, channel conversion is equivalent to mapping the distribution of latent sequence variables of the source sound $X^{(src)}$ to target sound $X^{(tar)}$. For each segment in $X^{(src)}$, we shift $Z_2^{(src,n)}$ by modifying μ_2 ($\Delta\mu_2 = \mu_2^{(tar)} - \mu_2^{(src)}$) while keeping $Z_1^{(src,n)}$ unaltered, as shown in Figure 4. Based on the method, we shift the z_2 values using the pre-obtained $\mu_2^{(tar)}$ (general mean of target domain DB) from the target domain, as shown in Figure 1.

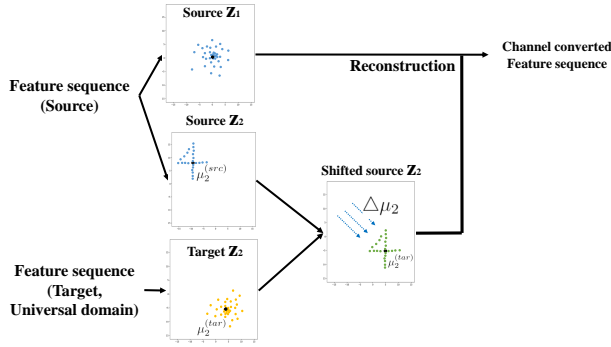


Fig. 4: Concept example of mean shifting for channel conversion

Here, as in the conventional domain adaptation [13], the target domain can be set to device A for channel conversion. (Device B, C to Device A) However, as mentioned above, input sounds might be unseen target domain input in the real-world environment, so modeling the relationship between specific domains cannot be a general solution for domain adaptation. Therefore, we propose an approach of converting unspecified domains into a target domain, not the approach of converting specific domains into a target domain. Following the proposed approach, the source domain (e.g. development-train set in DCASE 2018 task 1-B), which is generally a large size DB, also have to be converted. For the convenience of notation, we set the ‘universal domain’ to the our proposed target domain. As shown in Figure 5, we trained the FHVAE and obtained μ_2 of the universal domain through pre-training step. Using the μ_2 from the universal domain, channel conversion is conducted on the training DB and then the acoustic scene classifier is trained using the converted DB as shown in Figure 5. Since the channel converted DB is used for both training and testing step, the classifier is not affected by channel mismatching. Compare to the conventional method, the DB from the target domain is not required in the pre-training step for channel conversion.

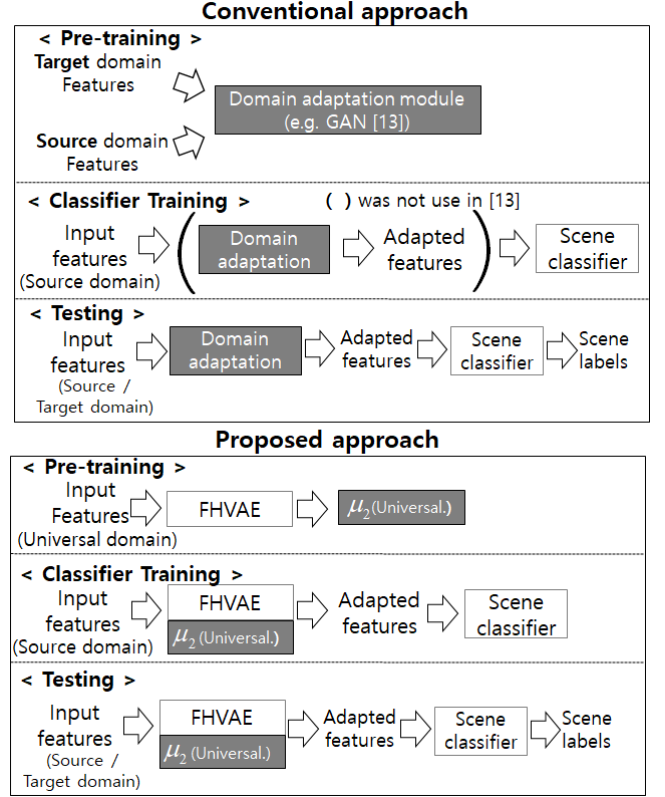


Fig. 5: Comparison of domain adaptation process between conventional and proposed approach

2.3. Acoustic scene classifier

For acoustic scene classification, we used the DCASE 2018 task 1 baseline system [9] to concentrate more on the performance changes of domain adaptation rather than classification performance itself. The baseline system implements a CNN based approach, where log mel scale-band energies are extracted for each 10-second signal, and a network consisting of two CNN layers and one fully connected layer is trained to assign scene labels to the audio signals. Detail configuration of baseline system can be found in the official website ¹.

3. EXPERIMENTAL SETTINGS AND RESULTS

For the acoustic scene classification experiment, DCASE 2018 task 1 dataset was used. The dataset was recorded in six large European cities, in different locations for each scene class. For each recording location, there are 5-6 minutes of audio. The original recordings were split into segments with a length of 10 seconds with one label from the pre-defined scene labels (Airport, Indoor shopping mall, Metro station, Pedestrian street, Public square, Street, Tram, Bus, Metro,

¹<http://dcase.community/challenge2018/task-acoustic-scene-classification>

Average accuracy on scene classes (%)	Device A (source)	Device B,C (target)	Averaged
DCASE 2018 baseline [9]	58.9	45.6	52.3
Domain adapted DCASE Kaggle model [13]	65.3	31.7	48.5
Proposed Methods			
Using Z_1 as feature (without reconstruction)	56.3	46.8	51.6
Shift z_2 to $\mu_2^{(tar)}$ of device A	58.2	47.1	52.7
Shift z_2 to $\mu_2^{(tar)}$ of device B and C	57.8	51.2	54.5
Shift z_2 to $\mu_2^{(tar)}$ of external scene DB (universal acoustic scene domain)	57.6	49.8	53.7

Table 1: Class average accuracy (%) on DCASE 2018 task 1-B development dataset. Note that the challenge was ranked by average accuracy on device B, C

and Urban park). The total size of the device A dataset is 8640 segments of 10 seconds length, i.e. 24 hours of audio. The dataset contains 2 hours of parallel data recorded with all three devices (A, B and C). The amount of data is as follows:

- Device A: 24 hours (8640 segments, 864 segments / scene)
- Device B: 2 hours (720 segments, 72 segments / scene)
- Device C: 2 hours (720 segments, 72 segments / scene)

The training subset contains 6122 segments from device A, 540 segments from device B, and 540 segments from device C. The test subset contains 2518 segments from device A, 180 segments from device B, and 180 segments from device C. Since the DB label set of the final evaluation DB set, which was used for challenge ranking, were not released, the training and testing were conducted with the DB configuration used during the challenge. We follow configurations of the FHVAE model in [17], but the input feature dimension is modified to the number of dimensions in the DCASE baseline ASC model. (log-mel of 40 dimensions) DB augmentation with SNR 10 and 15 dB level was conducted with white Gaussian noise in the FHVAE training. (3-times larger than original DB volume) We did not use augmented DB for classifier training. The experimental results are shown in Table 1.

For comparison with the baseline system, we first used z_1 latent variables, which showed channel invariant characteristics, as input features without reconstruction. System performance was slightly improved on device B and C, but performance for device A was worsened. This is because FHVAE extracts information about the residual generating factors that change from segment to segment for z_1 , rather than scene class discriminative information. On the other hand, since z_2 variables may have scene information, it is necessary to reconstruct the original scene information through the decoding process after converting the channel component as in the proposed method. The performance evaluation of the proposed method consists of three DB configurations. All the cases improved performance for device B and C, but performance was slightly decreased on device A set. First, similar to [13], we

set device A as a target domain and obtaining $\mu_2^{(tar)}$ for training and testing ASC system. Since the approach uses device A domain information, performance of device A set was not decreased much compared to baseline, but performance improvement of device B and C domain was limited. When device B and C domain information in development-train set were used in FHVAE, the highest ASC performance for device B, C was achieved. This is an obvious result because the system used the information of the target device in the pre-training step. In the last case, we obtained $\mu_2^{(tar)}$ using DB from DCASE 2016 and 2017 task 1 development set [7, 8]. Although the DB configuration and labels of DCASE 2016 and 2017 are different from the 2018 set, we could use these DBs, because the FHVAE system is based on unsupervised training. In this case, even though the FHVAE does not utilize domain information of the training and test DB of the scene classifier, a higher performance, compared to the case of using the device A set, was achieved using the various acoustic scene DBs. The results show the effectiveness of the proposed method, and it is noteworthy that the domain adaptation can be conducted only by using the universal domain information without the target or source domain information. Unlike the baseline model of DCASE 2018 task 1, a baseline model used in the previous domain adaptation research [13] more intensively trained on the device A domain rather than device B or C. Therefore, it is difficult to conduct a proper comparison experiment with the proposed methods. For comparison, it is necessary to conduct adversarial adaptation on the DCASE 2018 task 1 model [9] in the future work.

4. CONCLUSION AND FUTUREWORK

This paper proposes FHVAE based channel conversion for ASC in channel mismatching condition. Proposed latent variable shifting method shows performance improvement on the DCASE 2018 task 1-B DB. Especially, the proposed approach shows ASC performance improvement in mismatched condition by using only universal scene DB in pre-processing, without source domain (device A) and target domain (device B, C) information. For the future work, we plan to research for speech enhancement in a similar approach. In the situation of source domain (clean speech) and the target domain (unseen noisy speech), it would be meaningful if noisy channel adaptation is possible with only universal speech DBs, regardless of the target noise type.

Acknowledgements

The authors would like to thank Donmoon Lee and Dr. Yoonchang Han for valuable discussion to inspire ideas of the paper.

5. REFERENCES

- [1] Seongkyu Mun, Suwon Shon, Wooil Kim, David K. Han, and Hanseok Ko, “Deep neural network based learning and transferring mid-level audio features for acoustic scene classification,” in *ICASSP*, 2017, pp. 796–800.
- [2] Seongkyu Mun, Suwon Shon, Wooil Kim, and Hanseok Ko, “Deep Neural Network Bottleneck Features for Acoustic Event Recognition,” in *Interspeech*, 2016, pp. 2954–2957.
- [3] Wang Jun and Li Shengchen, “Self-Attention Mechanism Based System for Dcase2018 Challenge Task1 and Task4,” in *DCASE 2018 workshop*, 2018.
- [4] Zhao Ren and et. al., “Attention-Based Convolutional Neural Networks for Acoustic Scene Classification,” in *DCASE 2018 workshop*, 2018.
- [5] Seongkyu Mun, Sangwook Park, David K. Han, and Hanseok Ko, “Generative adversarial network based acoustic scene training set augmentation and selection using SVM hyper-plane,” in *DCASE Workshop*, 2017.
- [6] J.H Yang, Kim N.K, and Kim H.K, “Se-Resnet with Gan-Based Data Augmentation Applied to Acoustic Scene Classification,” in *DCASE 2018 workshop*, 2018.
- [7] Annamaria Mesaros, Toni Heittola, Emmanouil Benetos, Peter Foster, Mathieu Lagrange, Tuomas Virtanen, and Mark D. Plumbley, “Detection and Classification of Acoustic Scenes and Events: Outcome of the DCASE 2016 Challenge,” *IEEE/ACM Transactions on Audio Speech and Language Processing*, vol. 26, no. 2, pp. 379–393, 2018.
- [8] Annamaria Mesaros, Toni Heittola, and Tuomas Virtanen, “Acoustic scene classification: an overview of dcase 2017 challenge entries,” *International Workshop on Acoustic Signal Enhancement (IWAENC)*, 2018.
- [9] Annamaria Mesaros, Toni Heittola, and Tuomas Virtanen, “A multi-device dataset for urban acoustic scene classification,” in *DCASE Workshop*, 2018.
- [10] Hamid Eghbal-Zadeh, Bernhard Lehner, Matthias Dorfer, and Gerhard Widmer, “CP-JKU Submissions for DCASE-2016: a Hybrid Approach Using Binaural I-Vectors and Deep Convolutional Neural Networks,” in *DCASE Workshop*, 2016.
- [11] Yuma Sakashita and Masaki Aono, “Acoustic Scene Classification by Ensemble of Spectrograms Based on Adaptive Temporal Divisions,” Tech. Rep., DCASE 2018, 2018.
- [12] Yoonchang Han, Jeongsoo Park, and Kyogu Lee, “Convolutional neural networks with binaural representations and background subtraction for acoustic scene classification,” Tech. Rep., DCASE 2017, 2017.
- [13] Shayan Gharib, Konstantinos Drossos, Çakir Emre, Dmitriy Serdyuk, and Toumas Virtanen, “Unsupervised adversarial domain adaptation for acoustic scene classification,” in *DCASE Workshop*, 2018.
- [14] Wei-ning Hsu, Yu Zhang, and James Glass, “Unsupervised Learning of Disentangled and Interpretable Representations from Sequential Data,” in *Neural Information Processing Systems (NIPS)*, 2017.
- [15] Wei-ning Hsu, Yu Zhang, and James Glass, “Learning Latent Representations for Speech Generation and Transformation,” in *Interspeech*, 2017, pp. 1273–1277.
- [16] Wei-Ning Hsu and James Glass, “Extracting Domain Invariant Features by Unsupervised Learning for Robust Automatic Speech Recognition,” in *ICASSP*, 2018, pp. 5614–5618.
- [17] Wei-ning Hsu and James Glass, “Scalable Factorized Hierarchical Variational Autoencoder Training,” in *Interspeech*, 2018, pp. 1462–1466.
- [18] Diederik P. Kingma and Max Welling, “Auto-Encoding Variational Bayes,” in *International Conference on Learning Representations (ICLR)*, 2014, pp. 1–14.
- [19] Suwon Shon, Wei-Ning Hsu, and James Glass, “Unsupervised Representation Learning of Speech for Dialect Identification,” in *IEEE Spoken Language Technology Workshop (SLT)*, 2018.
- [20] Wei-ning Hsu, Hao Tang, and James Glass, “Unsupervised Adaptation with Interpretable Disentangled Representations for Distant Conversational Speech Recognition,” in *Interspeech*, 2018, pp. 1576–1580.