

# LEARNING A JOINT EMBEDDING SPACE OF MONOPHONIC AND MIXED MUSIC SIGNALS FOR SINGING VOICE

Kyungyun Lee      Juhan Nam

Graduate School of Culture Technology, KAIST  
{kyungyun.lee, juhannam}@kaist.ac.kr

## ABSTRACT

Previous approaches in singer identification have used one of monophonic vocal tracks or mixed tracks containing multiple instruments, leaving a semantic gap between these two domains of audio. In this paper, we present a system to learn a joint embedding space of monophonic and mixed tracks for singing voice. We use a metric learning method, which ensures that tracks from both domains of the same singer are mapped closer to each other than those of different singers. We train the system on a large synthetic dataset generated by music mashup to reflect real-world music recordings. Our approach opens up new possibilities for cross-domain tasks, e.g., given a monophonic track of a singer as a query, retrieving mixed tracks sung by the same singer from the database. Also, it requires no additional vocal enhancement steps such as source separation. We show the effectiveness of our system for singer identification and query-by-singer in both the in-domain and cross-domain tasks.

## 1. INTRODUCTION

Singing voice is often at the center of attention in popular music. We can easily observe large public interest in singing voice and singers through the popularity of karaoke industry and singing-oriented television shows. A recent study also showed that some of the most salient components of music are singers (vocals, voice) and lyrics [5]. Therefore, extracting information relevant to singing voice, i.e., to singers, from music signals, is an important area of research in music information retrieval (MIR) [9, 11]. The relevant tasks include singing voice detection [16], singing melody extraction [14, 28], singer identification [12, 21], and similarity-based music retrieval [8, 22].


Modern singer information processing systems have been designed to work with only one of monophonic or mixed music signals [15, 21, 34]. Then, given both types of signals for analysis, we question whether the system can extract information relevant to singing voice that is transferable between monophonic and mixed tracks. In

our experiment, we observe that systems trained with only one type of signals do not perform well, when tested with another type. To address this limitation, we introduce a new problem of *cross-domain* singer identification (singer-ID) and similarity-based retrieval, in which we regard monophonic and mixed music signals as two different audio domains. Cross-domain problems have been explored in computer vision and recommender systems, for example, image retrieval from user sketches to real images [29] and user preference modeling from movies to books [6]. In MIR, information transfer between monophonic and mixed tracks can open up new possibilities for singer-based retrieval systems. Some examples are: 1) given a user’s vocal recording in a karaoke application, finding popular singers who sound similar to the user, and 2) given a studio vocal track of a singer, retrieving all tracks (monophonic and mixed) relevant to the singer from a large music database.

To learn a joint feature representation of data from both monophonic and mixed tracks, we adopt a metric learning method, which forces tracks from the same singer to be mapped closer to each other than those from others (Section 3.2). To acquire sufficient training data, we create a synthetic dataset by performing a simple music mashup on two public datasets: vocal recordings from DAMP [30] and background tracks from musdb18 [25] (Section 3.1.1). We present experiments to demonstrate that our system is able to extract singer-relevant information from both monophonic and mixed music signals, and share the information between the two domains (Section 4). Source code, trained models, example audios and detailed information about the dataset are available<sup>1</sup>.

## 2. RELATED WORK

Cross-domain systems have not yet been examined regarding singing voice analysis. Nonetheless, a common challenge in singer information processing systems is to extract singing voice characteristics from music signals in the presence of background accompaniment music. The most direct way to obtain vocal information is to use monophonic vocal tracks. Recently, Wang et al. [34] trained a siamese neural network on monophonic recordings from a subset of the DAMP dataset. Their model scored higher on singer classification but lower on song classification compared to a baseline model. This implies that the model was able to learn singing voice characteristics, rather than the

 © Kyungyun Lee, Juhan Nam. Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). **Attribution:** Kyungyun Lee, Juhan Nam. “Learning a joint embedding space of monophonic and mixed music signals for singing voice”, 20th International Society for Music Information Retrieval Conference, Delft, The Netherlands, 2019.

<sup>1</sup> <http://github.com/kyungyunlee/mono2mixed-singer>

content of a music piece, such as its melody or lyrics. However, since music of our interest is often mixed tracks, this approach has limitations.

Several works have handled mixed audio signals by enhancing vocal signals through source separation or melody enhancement [8, 15, 21]. Given recent advances in source separation [32], this approach may bring improved results for most singing voice analysis systems. Another common choice is using audio features that represent human voice or singing voice, such as mel-frequency cepstral coefficients (MFCCs) [2, 15]. With the success of deep neural networks, it is even possible to learn appropriate features from more general audio representations, i.e., short-time Fourier transform (STFT) or even raw audio. We take this last approach and train our model to be a feature extractor for a given input audio represented by a mel-spectrogram.

Depending on the target task, background music can be helpful. An example is singer recognition in popular music [19]. This is because singing style is often dependent on the genre or mood of the music, and singers tend to perform in similar genres throughout their careers. However, our work focuses on learning the actual characteristics of singing voice, independent from background music.

### 3. METHODS

In this section, we describe the data generation pipeline, model configuration and training strategy for learning a joint representation of monophonic and mixed tracks for singing voice.

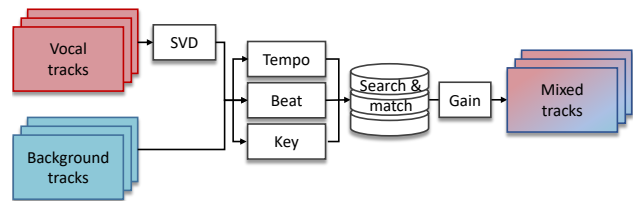
#### 3.1 Data generation

For training cross-domain singer-ID and retrieval systems, a sufficiently large number of monophonic and mixed track pairs per singer is needed. Existing singing voice datasets, such as MIR-1K [10], iKala [3] and Kara1K [1], provide the monophonic and mixed track pairs, but they have a small number of singers or only a few tracks per singer. An alternative option may be to perform singing voice detection (SVD) and vocal source separation on a large dataset, but the audio quality can be degraded.

In this work, we choose to utilize the DAMP dataset, which contains vocal-only recordings from mobile phones of around 3,500 users from the Smule karaoke app (there are 10 full-length songs per user). This serves as the main ingredient to generate our synthetic singer dataset. As a preprocessing step, we perform a simple energy-based SVD to remove silent segments. Then, 1000 singers are chosen for training stage and additional 300 singers are put aside for testing. The original DAMP dataset processed with SVD, *DAMP-Vocal*, is used as the monophonic dataset; the synthesized mixed track dataset, *DAMP-Mash* (detailed in section 3.1.1), is used as the mixed track dataset in this work.

##### 3.1.1 Mashup: DAMP and musdb18

A music mashup is a way of creating music by carefully mixing two or more tracks from several different pre-recorded songs. Inspired by such work, we automatically



**Figure 1:** Mashup pipeline to generate synthetic dataset, *DAMP-Mash*.

generate a synthetic singer dataset, called DAMP-Mash, by combining vocal recordings from the DAMP dataset with background tracks from the musdb18 dataset. Instead of random mixing, we build a pipeline (Figure 1) to identify the "mashability" [4] between tracks. Our mashability criteria requires 3-second long vocal and background tracks to have the same tempo and key. Once the two segments pass the mashability test, they are mixed at their nearest beat location. Before mixing, we adjust the loudness by balancing the root-mean-square energy between both segments.

Tempo detection and beat tracking are performed at track-level using librosa 0.6.2 [20]. On the other hand, key is determined locally at 3-second long segments by using the Krumhansl-Schmuckler key finding algorithm [13] on chromagrams. As a result, vocal segments within the same song end up being mixed with multiple different background tracks. Thus, we view our synthetic dataset as being genre-independent. This mashup pipeline can be also regarded as a data augmentation technique.

#### 3.2 Model

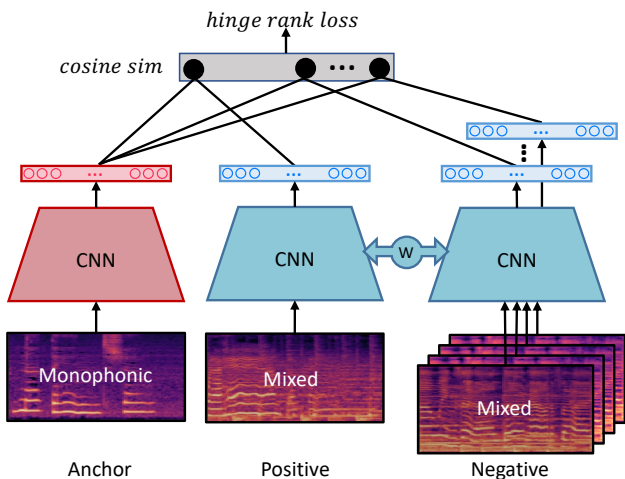
##### 3.2.1 Skeleton model

A 5-layer 1-D convolutional neural network (CNN) is the skeleton of larger networks used in this paper. First four convolutional layers have 128 filters of size 3, each followed by a maxpooling layer of size 3. The last convolutional layer consists of 256 filters of size 1 to output a final embedding vector of 256 dimensions. All convolution operations are done on the temporal dimension only. The final embedding vector will be used to represent input audio and to perform tasks described in Section 4. Batch normalization and Leaky ReLU [17] are applied to all convolutional layers, and dropout of 50% is applied after the last convolutional layer to prevent overfitting.

The input is a 3-second long audio of at least 70% singing voice frames. The sample rate of audio is 22050 Hz and we compute an STFT using a 1024-sample long Hanning window with 50% overlap. We then convert it to a mel-spectrogram with 128 bins and apply logarithmic compression on the magnitude. As a result, the input shape is 129 frames by 128 bins.

##### 3.2.2 Embedding model

The outcome of metric learning is a mapping function from inputs to output vectors in an embedding space, where inputs of the same class are closer to each other than those



**Figure 2:** Configuration of CROSS model. The anchor network (left) is for modeling monophonic tracks; the rest is for modeling mixed tracks.

from different classes. Specifically, we build our model upon a triplet network, which consists of three potentially weight-sharing networks that takes three inputs: anchor, positive (same class as the anchor) and negative (different class as the anchor) items. This architecture can be extended to take multiple negative items [31] to overcome the limitation of learning from only one negative example. For model configuration, we closely follow the work of [24], using 4 negative items. We also use a type of margin-based ranking loss, called hinge rank loss, with cosine similarity as our metric [7].

Thus, the loss function for a given set of anchor ( $p_i$ ), positive ( $p_+$ ) and negative ( $p_-$ ) feature vectors is:

$$loss(p_i, p) = \sum_{p_-} \max[0, \alpha - S(p_i, p_+) + S(p_i, p_-)] \quad (1)$$

where  $S$  is a similarity score:

$$S(p_i, p_j) = \cos(p_i, p_j) = \frac{p_i \cdot p_j}{\|p_i\| \cdot \|p_j\|} \quad (2)$$

and  $\alpha$  indicates the margin, which is fixed to 0.1 after performing a grid search on values between 0.01 and 1.0. Negative tracks are selected through negative sampling among tracks that do not belong to the singer of the anchor item. We tested a more difficult negative sampling strategy of selecting the four highest scoring items among twenty randomly chosen negative samples, but the model showed minor improvement with an increase in computation time. Investigation on negative sampling is left as our future work.

We choose metric learning for three main reasons. First, by giving a higher similarity score to any pair of tracks performed by the same singer, the model can learn to identify the singer from a track regardless of it being monophonic or mixed. Thus, it is especially suitable for training cross-domain systems. Second, using singer identity as the only ground truth to measure similarity between two tracks will force the model to focus only on singing voice. Since DAMP-Mash is genre-invariant, the only common component in two tracks is going to be related to singing voice.

Thus, we may expect the model to perform a feature-level source separation on mixed tracks. Lastly, the model can be trained on a larger number of singers without increasing the number of parameters. On the other hand, a classification model that uses a softmax layer will need to increase the output layer size to match the number of training singers [24].

We explore our ideas with three models, which differ in the type of data used for training:

- MONO: all inputs are monophonic tracks
- MIXED: all inputs are mixed tracks
- CROSS: anchors are monophonic, while positive and negative items are mixed tracks (Figure 2)

Our main idea in this work is reflected in the CROSS model, for which the hinge rank loss ensures that the cosine similarity between monophonic and mixed tracks from the same singer is scored higher than that from different singers. MONO and MIXED models are reference models for comparison.

While networks within MONO and MIXED model share weights, in CROSS model, the anchor network and the rest do not share weights. Thus, it yields two separate feed-forward networks, each designed specifically for its corresponding domain (Figure 2). As a result, depending on the domain of an input audio, one of the two networks is used as the feature extractor. Each network is configured with the skeleton model described in Section 3.2.1.

### 3.2.3 Pre-training via classification

Metric learning is known for its difficulty in optimization [35, 36]. To alleviate this problem, we train a classification model and use it to initiate the learning of the embedding models. The classification model has one linear layer added to the skeleton model (Section 3.2.1) and predicts the correct singer with a softmax probability. Instead of fully training it, we remove the last output layer after 30 epochs and use it to continue the training in a metric learning style. We do not freeze any layers.

## 4. EXPERIMENTS & EVALUATION

### 4.1 Test scenarios

Two main tasks for evaluation are singer identification and query-by-singer. In both tasks, a music signal to be analyzed (source) is queried to a collection of data (target) to retrieve desired information. Depending on the domain of source and target data, we design three test scenarios:

- *Mono2Mono*: both source and target data are *monophonic* (in-domain)
- *Mix2Mix*: both source and target data are *mixed* (in-domain)
- *Mono2Mix*: source data is *monophonic*, but the target data is *mixed* (cross-domain)

Each task is evaluated on all three test scenarios.

## 4.2 Task 1: Singer identification

**Dataset** : We select 300 singers unseen from the training stage for evaluation. For each singer, we use 6 tracks for building singer models and set aside 4 tracks as query tracks, resulting in 1200 queries. Depending on the domain of source and target data, DAMP-Vocal (monophonic) and DAMP-Mash (mixed) dataset are used accordingly.

**Description** : As in [23, 27], singer identification is to determine the correct singer of the query track among the 300 candidate singer models. All queries and singer models are represented as 256 dimensional feature vectors; a track vector is an average of 20 feature vectors computed from 3-second long segment of the same track and a singer model is an average of 6 track vectors from the same singer. We made predictions by computing cosine similarity (2) between the query track and all the singer models. Then, the singer with the highest score is chosen.

For our baseline, we train a Gaussian Mixture Model-Universal Background Model (GMM-UBM), which is commonly used in speaker recognition systems [26]. Each singer model is adapted through maximum a posteriori (MAP) estimation from a single singer-independent background model. All models are composed of 256 components with MFCCs of dimension 13 as input. We train 2 GMM-UBMs, one with monophonic tracks for *Mono2Mono* and the other with mixed tracks for *Mix2Mix*.

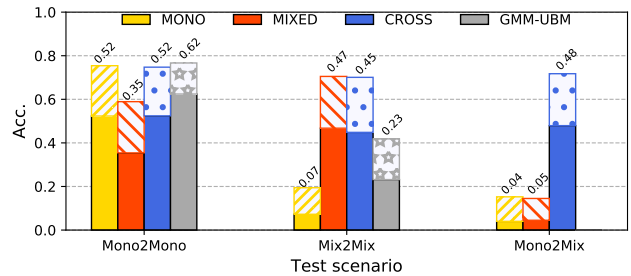
We report both top-1 and top-5 classification accuracy. They represent the proportion of correct guesses out of 1200 queries in total. Top-5 accuracy is calculated by considering a prediction as being correct if the ground truth singer appears within the top 5 highest scoring singers.

## 4.3 Task 2: Query-by-singer

**Dataset** : As in the singer recognition task (Section 4.2), same 300 singers are used for evaluation. 6 tracks from each singer are selected to build a collection of 1800 tracks to represent a search database and 4 tracks are used as test queries.

**Description** : Given a query track, the task is to retrieve tracks that are performed by the same singer among the track database. We compute the similarity (Equation (2)) between the query track and all the tracks in the database, and rank them based on their similarity scores. This can be applied to singer-based music recommender systems to discover singers with similar singing voice characteristics.

We report *precision* and *recall-at-k* as well as mean average precision (mAP) score, where  $k$  is set to 5 to resemble music retrieval systems. Given a query track performed by singer  $A$ , *precision-at-k* ( $Pr@k$ ) refers to the proportion of tracks that are performed by  $A$  and are recommended among  $k$  items; *recall-at-k* ( $R@k$ ) refers to the proportion of tracks that are performed by  $A$  and are recommended out of all the tracks performed by  $A$  (6 tracks in total) in the database. Unlike  $Pr@k$  and  $R@k$ , mAP takes into account the actual order of the recommended tracks. Thus, it



**Figure 3:** Singer identification results for MONO, MIXED and CROSS on different test scenarios. The solid section points to the top-1 accuracy (also written above each bar) and the hatched section points to the top-5 accuracy.

is useful for music recommender systems, where it is important that relevant items are not only retrieved, but also with higher confidence than false positive items.

## 4.4 Results

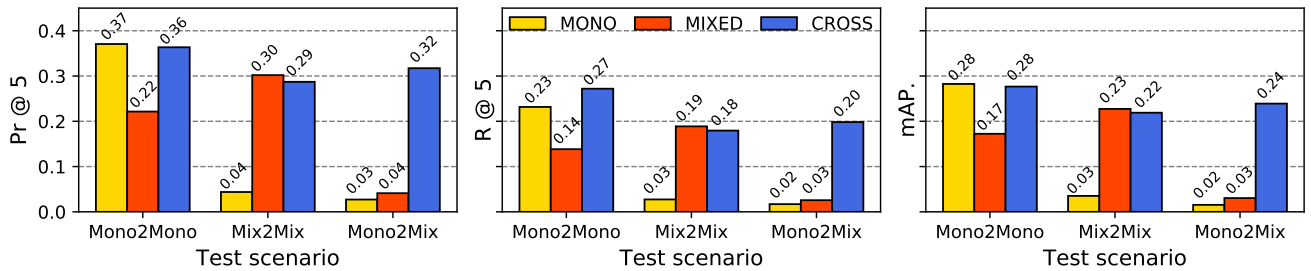
In Figure 3 and Figure 4, we observe a large performance variation for MONO and MIXED models across different test scenarios. Both models perform well in *Mono2Mono* and *Mix2Mix*, respectively. However, their performances drop significantly in other scenarios, especially for *Mono2Mix*. This is expected, because these models have not been trained to handle cross-domain scenarios.

On the other hand, CROSS model performs well on all three test scenarios, benefiting from two jointly trained networks that can each handle monophonic and mixed tracks. We see that it is the only model that is able to match and compare the singer identity between tracks from different domains. Also, its performances on *Mono2Mono* and *Mix2Mix* are on par with the MONO and MIXED models. This is a useful observation, since training only the CROSS model can still give good performance on all three test scenarios, avoiding the effort of training separate models for each scenario. Note that the baseline model, GMM-UBM, shows the best performance in *Mono2Mono*, but not so well in *Mix2Mix*. Result for *Mono2Mix* is omitted, since it is close to random prediction. When there is no background music, GMM-UBM with MFCCs are efficient in characterizing singing voice.

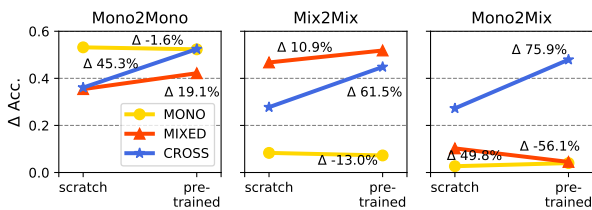
As mentioned in Section 3.2.3, we show the effect of using a pre-trained network on singer identification task (Figure 5). CROSS model (blue stars) shows the largest performance improvement compared to the other two models. We assume that comparing the singer identity between monophonic and mixed track is more difficult than comparing between tracks of the same domain. Therefore, a pre-trained model, which learned to somewhat identify singers from an input audio, serves as a hint to focus on signals relevant to singing voice. Using a pre-trained model not only improved the accuracy, but also accelerated the learning process.

Regarding background music as noise and singing voice as the signal, signal-to-noise ratio (SNR) has a large impact on the performance of singing voice analysis systems [16]. We change the SNR of the test data and show

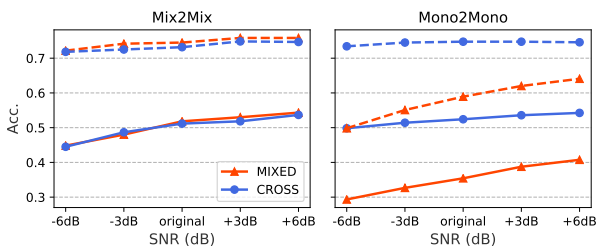




**Figure 4:** Query-by-singer result for MONO, MIXED and CROSS models. *precision-at-k* (left), *recall-at-k* (center) and mean average precision score(right) are shown. Each number above the bar refers to corresponding model’s score.



**Figure 5:** Top-1 accuracy improvement on singer identification when using a pre-trained network. The numbers on the line indicates the percentage of improvement.



**Figure 6:** Result of MIXED and CROSS model for singer recognition on varying SNR. Solid line indicates top-1 and dotted line indicates top-5 accuracy.

results on singer recognition task for MIXED and CROSS models in Figure 6. Since *Mono2Mono* deals with only monophonic tracks, the change in performance exhibited in *Mono2Mono* (right) is due to the overall loudness of the track, not SNR. Therefore, as the performance change on *Mix2Mix* (left) shows a similar trend across different SNR, it implies that models trained on DAMP-Mash dataset is able to identify singing voice in more difficult conditions. This is a great benefit for singing voice analysis systems.

#### 4.5 Evaluation on Popular Music Recordings

As our system is trained with a synthetic dataset, we evaluate it on popular music recordings to ensure that the trained system can also generalize to real-world data.

**Dataset :** Million Song Dataset (MSD) contains 1,000,000 tracks and 44,745 artists from popular music recordings. As done in [23], we filter the dataset to select artists with substantial vocal tracks using singing voice detec-

	MIXED	CROSS	POP	Random
Acc.	0.291	0.282	0.393	0.002
Top-5 Acc.	0.511	0.491	0.664	0.01

**Table 1:** Accuracy result on singer recognition on dataset of popular music recording, MSD-Singerdataset

tion (SVD). This dataset is referred to as *MSD-Singer*<sup>2</sup>. For comparison, we train a model, named POP, on 1000 artists from MSD-Singer dataset. We used 17 30-second long tracks for each artist for training. 500 singers, unseen from the training stage, are used for evaluation. 15 tracks of each singer are used for building singer models and 5 tracks are used as query tracks.

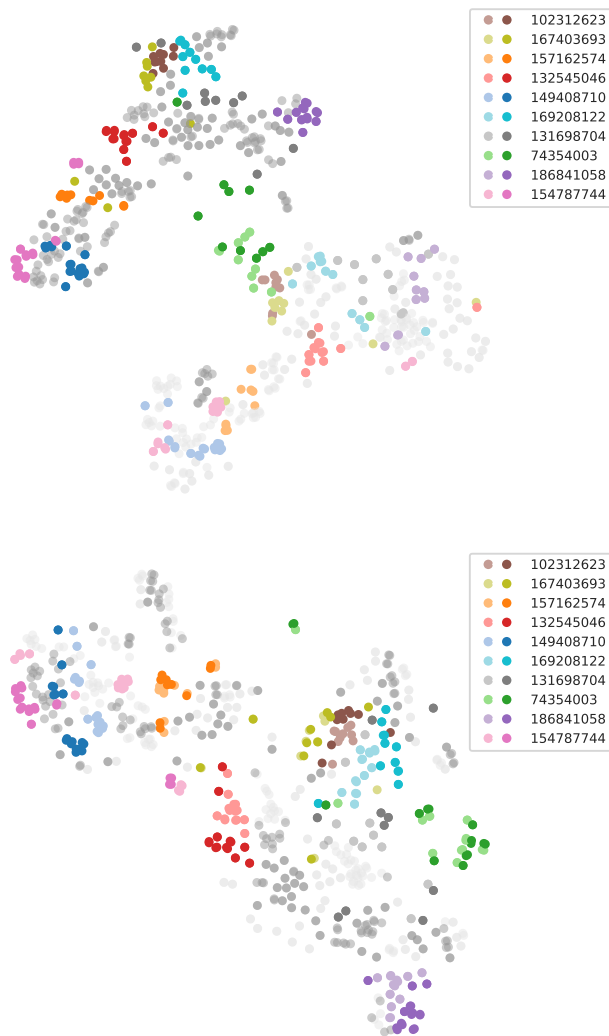
**Description :** The task is equivalent to singer identification in Section 4.2, only with a different dataset. The result from the POP model is the upper bound, as it is trained and tested on MSD-Singer dataset; meanwhile, MIXED and CROSS models are trained on DAMP-Mash and DAMP-Vocal dataset.

**Result:** The result shown in Table 1 compares MIXED and CROSS models with POP model and a random classifier. It shows that even though our models are trained only with the synthetic dataset, they are also able to identify singing voice in popular music. Therefore, we can confirm that DAMP-Mash dataset is able to represent the popular music to some degree and that our models are able to generalize to real-world recordings. We believe that the results will improve with a better automatic mashup pipeline. Training the CROSS model on source separated MSD-Singer dataset is also left as a future work.

## 5. EMBEDDING SPACE VISUALIZATION

We visualize the embedding space learned by the MIXED and CROSS models to understand how they each process monophonic and mixed tracks. From DAMP-Voice and DAMP-Mash dataset, we select 25 singers unseen from the training stage and highlight 10 with colors for better visualization. 20 tracks are plotted for each singer: 10 monophonic vocal tracks and their corresponding mixed tracks. After feature extraction, we reduce the dimension of each

<sup>2</sup> Details provided at <http://github.com/kyungyunlee/MSD-singer>



**Figure 7:** Singer embedding space from MIXED model (top) and CROSS model (bottom). The label numbers are player IDs from the DAMP dataset. The colors on the left column refers to monophonic vocal tracks; the right column refers to mixed tracks. Best viewed in color.

feature vector from 256 to 2 dimensions using t-distributed stochastic neighbor embedding (t-SNE) [18]. Each dot on the embedding space represents a track. For visualization, we use a paired color palette and assign lighter color to monophonic tracks. Since both monophonic and mixed tracks are from the same singer, an ideal embedding space will show clusters of 20 tracks for each singer.

Figure 7 shows the embedding space learned from the MIXED model. There is a noticeable gap between the features of monophonic tracks and that of mixed tracks, which means that the model differentiates monophonic and mixed tracks, rather than finding similar singing voice. Still, we can see that the model is able to cluster tracks from the same singer within the same domain. However, in Figure 7, the monophonic and mixed track features of the same singer are mapped close to each other. This explains why the CROSS model shows good performance on cross-domain tasks. We can observe that it is able to trans-

fer singer information across two domains.

## 6. MOTIVATION FOR FUTURE WORK

**Improvement on music mashup:** Our mashup pipeline has a large room for improvement. Besides errors produced from existing algorithms, such as key detection, more efforts can be put towards mixing two tracks with a good balance as in real-world recordings. A good automatic mashup system can benefit many areas of research in MIR. The creativity and limitless choices of techniques that can be applied to generate a mashup imply that a large amount of multitrack dataset can be generated for many tasks of interest.

**From track to singer modeling:** In this work, we use an average of several track-level feature vectors to build singer models. However, in case of singers with highly varying vocal characteristic between different tracks and taking into account the “album effect”, averaging may not always be the best choice. Exploring GMMs with multiple mixtures or principal component analysis (PCA) can be an interesting future direction.

**Going beyond singing voice:** Although we have focused on singing voice, our methods can be tested with tasks involving other instruments, such as multiple instrument recognition. The same mashup technique can be applied to create a dataset, by replacing the monophonic vocal tracks with any instrument of interest. Data generation with mashup may yield better results for instrument recognition in real-world recordings compared to the method proposed in [33], where only two monophonic instrument tracks are used to create a random mix.

## 7. CONCLUSION

In this paper, we introduced a new problem of cross-domain singer identification and singer-based music retrieval to allow information transfer between monophonic and mixed tracks. Through data generation using music mashup, we were able to train an embedding model to output a joint representation for singing voice from tracks regardless of their domain. We evaluated on three different test scenarios, which include both in-domain and cross-domain cases. A huge advantage of CROSS model is that it performs well not only on the cross-domain scenario, but also on commonly observed in-domain scenarios. Therefore, by training only the CROSS model, it yields two models, one for each domain. Additional evaluation on varying SNR and on popular music dataset demonstrated that the model is robust to background music and can also be generalized beyond our synthetic dataset.

To conclude, we believe that cross-domain systems can enable many interesting applications related to singing voice, as well as in MIR. Specifically, our future interests include improving the quality of the mashup dataset and performing comparisons between singing voices of karaoke users and that of popular singers.

## 8. ACKNOWLEDGEMENTS

We thank Keunwoo Choi for valuable comments and reviews. This work was supported by Basic Science Research Program through the National Research Foundation of Korea funded by the Ministry of Science, ICT & Future Planning (2015R1C1A1A02036962), and by NAVER Corp.

## 9. REFERENCES

- [1] Yann Bayle, Ladislav Maršík, Martin Rusek, Matthias Robine, Pierre Hanna, Katerina Slaninová, Jan Marti-novic, and Jaroslav Pokorný. Kar1k: a karaoke dataset for cover song identification and singing voice analy-sis. In *IEEE International Symposium on Multimedia (ISM)*, pages 177–184, 2017.
- [2] Wei Cai, Qiang Li, and Xin Guan. Automatic singer identification based on auditory features. In *2011 Sev-enth International Conference on Natural Computa-tion*, volume 3, pages 1624–1628, 2011.
- [3] Tak-Shing Chan, Tzu-Chun Yeh, Zhe-Cheng Fan, Hung-Wei Chen, Li Su, Yi-Hsuan Yang, and Roger Jang. Vocal activity informed singing voice separation with the iKALA dataset. In *IEEE International Con-ference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 718–722. IEEE, 2015.
- [4] Matthew EP Davies, Philippe Hamel, Kazuyoshi Yoshii, and Masataka Goto. Automashupper: Au-tomatic creation of multi-song music mashups. *IEEE/ACM Transactions on Audio, Speech, and Lan-guage Processing*, 22(12):1726–1737, 2014.
- [5] Andrew Demetriou, Andreas Jansson, Aparna Kumar, and R Bittner. Vocals in music matter: The relevance of vocals in the minds of listeners. In *Proceedings of the 19th International Society for Music Information Retrieval Conference (ISMIR)*, pages 514–520, 2018.
- [6] Ali Mamdouh Elkahky, Yang Song, and Xiaodong He. A multi-view deep learning approach for cross do-main user modeling in recommendation systems. In *Proceedings of the 24th International Conference on World Wide Web*, pages 278–288, 2015.
- [7] Andrea Frome, Greg S Corrado, Jon Shlens, Samy Bengio, Jeff Dean, Tomas Mikolov, et al. Devise: A deep visual-semantic embedding model. In *Ad-vances in Neural Information Processing Systems*, pages 2121–2129, 2013.
- [8] Hiromasa Fujihara, Tetsuro Kitahara, Masataka Goto, Kazunori Komatani, Tetsuya Ogata, and Hiroshi G Okuno. Singer identification based on accompaniment sound reduction and reliable frame selection. In *Pro-ceedings of the 6th International Conference on Music Information Retrieval (ISMIR)*, pages 329–336, 2005.
- [9] Masataka Goto. Singing information processing. In *Proceedings of the 12th IEEE International Confer-ence on Signal Processing (ICSP)*, volume 10, pages 2431–2438, 2014.
- [10] Chao-Ling Hsu and Jyh-Shing Roger Jang. On the improvement of singing voice separation for monau-ral recordings using the mir-1k dataset. *IEEE Trans-actions on Audio, Speech, and Language Processing*, 18(2):310–319, 2010.
- [11] Eric J Humphrey, Sravana Reddy, Prem Seetharaman, Aparna Kumar, Rachel M Bittner, Andrew Demetriou, Sankalp Gulati, Andreas Jansson, Tristan Jehan, Bern-hard Lehner, et al. An introduction to signal processing for singing-voice analysis: High notes in the effort to automate the understanding of vocals in music. *IEEE Signal Processing Magazine*, 36(1):82–94, 2019.
- [12] Youngmu Kim and Brian Whitman. Singer identifica-tion in popular music recordings using voice coding features. In *Proceedings of the 3rd International Con-ference on Music Information Retrieval*, 2002.
- [13] Carol L. Krumhansl. *Cognitive Foundations of Musi-cal Pitch*. Oxford psychology series. Oxford University Press, USA, 1990.
- [14] Sangeun Kum and Juhan Nam. Joint detection and classification of singing voice melody using convo-lutional recurrent neural networks. *Applied Sciences*, 9(7), 2019.
- [15] Mathieu Lagrange, Alexey Ozerov, and Emmanuel Vincent. Robust singer identification in polyphonic music using melody enhancement and uncertainty-based learning. In *Proceedings of the 13th Interna-tional Society for Music Information Retrieval Confer-ence (ISMIR)*, 2012.
- [16] Kyungyun Lee, Keunwoo Choi, and Juhan Nam. Re-visiting singing voice detection: a quantitative review and the future outlook. In *Proceedings of the 19th International Society for Music Information Retrieval Conference (ISMIR)*, 2018.
- [17] Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acous-tic models. In *Proceedings of the 30th International Conference on Machine Learning (ICML)*, 2013.
- [18] Laurens van der Maaten and Geoffrey Hinton. Visu-alizing data using t-sne. *Journal of Machine Learning Research*, 9:2579–2605, 2008.
- [19] Namunu Chinthaka Maddage, Changsheng Xu, and Ye Wang. Singer identification based on vocal and in-strumental models. In *Proceedings of the 17th Interna-tional Conference on Pattern Recognition (ICPR)*, vol-ume 2, pages 375–378, 2004.

- [20] Brian McFee, Matt McVicar, Stefan Balke, Carl Thomé, Vincent Lostanlen, Colin Raffel, Dana Lee, Oriol Nieto, Eric Battenberg, Dan Ellis, Ryuichi Yamamoto, Josh Moore, WZY, Rachel Bittner, Keunwoo Choi, Pius Friesch, Fabian-Robert Stöter, Matt Vollrath, Siddhartha Kumar, nehz, Simon Waloschek, Seth, Rimvydas Naktinis, Douglas Repetto, Curtis "Fjord" Hawthorne, CJ Carr, João Felipe Santos, JackieWu, Erik, and Adrian Holovaty. *librosa/librosa: 0.6.2*, August 2018.
- [21] Annamaria Mesaros, Tuomas Virtanen, and Anssi Klauri. Singer identification in polyphonic music using vocal separation and pattern recognition methods. In *Proceedings of the 8th International Conference on Music Information Retrieval (ISMIR)*, pages 375–378, 2007.
- [22] Tomoyasu Nakano, Kazuyoshi Yoshii, and Masataka Goto. Vocal timbre analysis using latent dirichlet allocation and cross-gender vocal timbre similarity. In *Acoustics, Speech and Signal Processing, 2014. ICASSP 2014. IEEE International Conference on*, pages 5202–5206, 2014.
- [23] Jiyoung Park, Donghyun Kim, Jongpil Lee, Sangeun Kum, and Juhan Nam. A hybrid of deep audio feature and i-vector for artist recognition. In *Joint Workshop on Machine Learning for Music, International Conference on Machine Learning*, 2018.
- [24] Jiyoung Park, Jongpil Lee, Jangyeon Park, Jung-Woo Ha, and Juhan Nam. Representation learning of music using artist labels. In *Proceedings of the 19th International Society for Music Information Retrieval Conference (ISMIR)*, 2018.
- [25] Zafar Rafii, Antoine Liutkus, Fabian-Robert Stöter, Stylianos Ioannis Mimilakis, and Rachel Bittner. The MUSDB18 corpus for music separation, December 2017.
- [26] Douglas A Reynolds, Thomas F Quatieri, and Robert B Dunn. Speaker verification using adapted gaussian mixture models. *Digital signal processing*, 10(1-3):19–41, 2000.
- [27] Jimena Royo-Letelier, Romain Hennequin, Viet-Anh Tran, and Manuel Moussallam. Disambiguating music artists at scale with audio metric learning. In *Proceedings of the 19th International Society for Music Information Retrieval Conference (ISMIR)*, Paris, France, 2018.
- [28] Justin Salamon, Emilia Gómez, Daniel PW Ellis, and Gaël Richard. Melody extraction from polyphonic music signals: Approaches, applications, and challenges. *IEEE Signal Processing Magazine*, 31(2):118–134, 2014.
- [29] Abhinav Shrivastava, Tomasz Malisiewicz, Abhinav Gupta, and Alexei A Efros. Data-driven visual similarity for cross-domain image matching. In *ACM Transactions on Graphics (ToG)*, volume 30, page 154, 2011.
- [30] Jeffrey C Smith. Correlation analyses of encoded music performance. 2013.
- [31] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *Advances in Neural Information Processing Systems*, pages 1857–1865, 2016.
- [32] Daniel Stoller, Sebastian Ewert, and Simon Dixon. Wave-u-net: A multi-scale neural network for end-to-end audio source separation. In *Proceedings of the 19th International Society for Music Information Retrieval Conference (ISMIR)*, Paris, France, 2018.
- [33] Yuji Tokozume, Yoshitaka Ushiku, and Tatsuya Harada. Learning from between-class examples for deep sound recognition. In *International Conference on Learning Representations (ICLR)*, 2018.
- [34] Cheng-i Wang and George Tzanetakis. Singing style investigation by residual siamese convolutional neural networks. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 116–120, 2018.
- [35] Jian Wang, Feng Zhou, Shilei Wen, Xiao Liu, and Yuanqing Lin. Deep metric learning with angular loss. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2593–2601, 2017.
- [36] Xiaofan Zhang, Feng Zhou, Yuanqing Lin, and Shaoting Zhang. Embedding label structures for fine-grained feature representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1114–1123, 2016.