

LOW-RESOURCE MUSIC GENRE CLASSIFICATION WITH CROSS-MODAL NEURAL MODEL REPROGRAMMING

Yun-Ning Hung¹, Chao-Han Huck Yang¹, Pin-Yu Chen², Alexander Lerch¹

¹ Georgia Institute of Technology, Atlanta, GA, USA

²IBM Research, Yorktown Heights, NY, USA

ABSTRACT

Transfer learning (TL) approaches have shown promising results when handling tasks with limited training data. However, considerable memory and computational resources are often required for fine-tuning pre-trained neural networks with target domain data. In this work, we introduce a novel method for leveraging pre-trained speech models for low-resource music classification based on the concept of *Neural Model Reprogramming (NMR)*. NMR aims at re-purposing a pre-trained model from a source domain to a target domain by modifying the input of a frozen pre-trained models for cross-modal adaptation. In addition to the known, input-independent, re-programming method, we propose a new reprogramming paradigm: *Input-dependent NMR*, to increase adaptability to complex input data such as musical audio. Experimental results suggest that a neural model pre-trained on large-scale datasets can successfully perform music genre classification by using this reprogramming method. The two proposed Input-dependent NMR TL methods outperform fine-tuning-based TL methods on a small genre classification dataset.

1. INTRODUCTION

Large-scale datasets are often recognized as one key component in successfully building powerful deep neural network (DNN) prediction models [1]. For example, the ImageNet dataset [2] with 14 million samples can be used to train several benchmark image classification systems [3] for visual perception. For sequence modeling and language processing tasks, the representation power of models pre-trained on millions of language data such as BERT [4] demonstrates promising results on few-shot tasks and low-resource prediction. In speech and acoustic understanding applications, the performance of acoustic models also benefits considerably from large-scale datasets such as human speech commands [5] and audio event prediction datasets (e.g., AudioSet [6]).

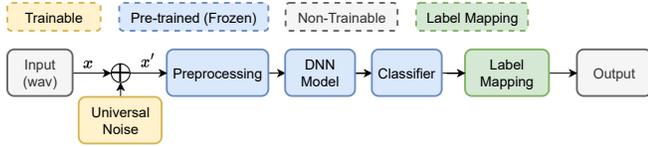
In the music domain, however, the lack of large-scale training data has been a critical problem. This not only hinders the development of data-driven approaches, but also makes the investigation of novel deep learning architectures (e.g., transformer-based models) developed in other domains challenging since they often require massive amounts of training

data. To tackle this problem, Transfer Learning (TL) is a popular solution. For example, the MusiCNN [7] and JukeBox [8] models, pre-trained large-scale music datasets, have achieved promising performance on several low-resource MIR downstream tasks [9]. VGGish and the OpenL3, pre-trained on the large-scale audio dataset, AudioSet [10, 11], also have shown their effectiveness on various music information retrieval (MIR) downstream tasks [12, 13].

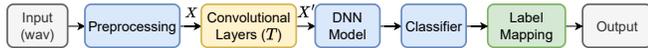
Although these TL methods could demonstrate decent performance in downstream tasks, they suffer from some drawbacks. For example, the learned representation does not guarantee to contain task-specific information. Recent results show, e.g., that models pre-trained on music auto-tagging might lack of information for key detection [14]. Moreover, it seems that learned representations might still lead to inferior performance compared to task-specific models designed for the purpose [14]. Theoretically, fine-tuning the pre-trained model should allow to solve this problem. However, training large-scale models such as JukeBox or VGGish requires considerable computational resources.

In this work, we attempt to tackle this problem from a different angle. We take advantage of a newly proposed method called Neural Model Reprogramming (NMR) [15, 16, 17], to adopt a pre-trained model for downstream tasks. NMR is an alternative TL technique that has been confirmed to provide good or even state-of-the-art results [16, 17] in numerous machine learning tasks. It aims to re-purpose a (frozen) pre-trained model on a task-specific dataset with only a small amount of training parameters without modifying the whole model. Unlike the original NMR training scheme, which proposes adding an universal trainable noise directly to the input sequence, we further propose two novel input-dependent methods to learn sample-dependent information from different input signals and improve the NMR training strategy on musical data.

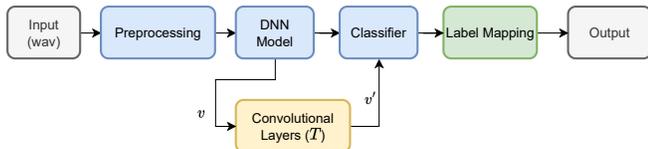
We investigate the proposed method on two models pre-trained on audio and speech, respectively: Audio Spectrogram Transformer (AST) [18] and SpeechATT [19], for the task of music genre classification. By using the proposed methods, we achieve new state-of-the-art results on the small music genre classification dataset, GTZAN [20].



(a) Training pipeline for the baseline input-independent NMR by adding an universal noise to the input waveform (II-NMR).



(b) Training pipeline for the proposed input-dependent NMR by transforming the input feature (ID-NMR).



(c) Training pipeline for the proposed input-dependent NMR method with skip connection (IDS-NMR).

Fig. 1: The overview of the training pipeline for the three NMR methods used in this work.

2. RELATED WORK

The concept of NMR is originally related to adversarial attacks [21]. By applying “trainable perturbations” against the loss function of a target model on the input data, an adversarial noise can maliciously manipulate output predictions of neural network based classification model. Meanwhile, Elsayed et al. [15] proposed using this trainable noise to *reprogram* a model pre-trained on a specific “source” task (e.g., ImageNet) to a new “target” task (e.g., MNIST or CIFAR-10). After the reprogramming model is successfully trained on the target data, these reprogramming noises can be directly applied to a frozen pre-trained model input during inference time. These trainable perturbations can thus be considered a new *program* to empower cross-domain transfer learning toward new tasks. NMR also demonstrates its flexibility in an application-oriented setting [22], such as reprogramming a frozen traffic sign image prediction API into a medical image classifier without knowing the parameters of a source model.

In addition to image data, the NMR method has been proven effective in various domains. In natural language processing, NMR has been proposed for machine translation and sentiment classification [23]. Vinod et al. further explore the possibility of cross-domain reprogramming by adopting NLP models to molecule learning tasks in biochemistry [17]. In terms of time series processing, the recent advances [16, 25] of Voice2Series show an English speech command model could be re-programmed to become either a sensor data predictor or a multilingual recognizer on some low-resource language such as Lithuanian with state-of-the-art performance.

The aforementioned techniques provide both preliminary empirical findings and theoretical foundations to further motivate us to study NMR techniques toward music applications, addressing the limitations of insufficient training data and infeasible large-scale neural architectures.

3. NEURAL MODEL REPROGRAMMING

This work introduces two new reprogramming methods. While traditional reprogramming simply adds noise to the input of a system as described in Sect. 3.1, our methods transform the input signal non-linearly and utilize intermediate representations to account for the variability and complexity of musical input signals. As shown in Fig. 1, all our pre-trained systems include input waveform, preprocessing step to extract a high-level representation, deep neural network, and a classifier to predict the probability of each label. The output of the classifier is then mapped by a many-to-one label mapping layer to map the probabilities of n source classes to a target class. That is, for a target label $y_T \in Y_T$, the prediction will be the averaged class predictions $(\{y_1, y_2, \dots, y_n\} \in Y_S)$ over the set of source labels from the original pre-trained model assigned to it. All of the steps mentioned above will not be updated during training. Instead, a trainable transformation function \mathcal{H} will be added to introduce “noise” to the training pipeline. We will introduce how to add \mathcal{H} .

3.1. Input-independent NMR: Baseline Method

Since the applicability of NMR to music signals remains unexplored, we first investigate the efficiency and potential of NMR for music signals by establishing a baseline. We first follow the same training scheme proposed by Yang et al. [16] as one simple baseline of waveform-level reprogramming. As shown in Fig. 1a, a trainable universal parameter, θ , is added as a “noise” to transform the whole waveform: $x' = \mathcal{H}(x; \theta) := x + \theta$, where x and x' represent the waveform before and after transformation, respectively. Since the universal parameter is independent of the input, we refer to this method as *II-NMR*.

3.2. Input-dependent NMR

Although most previous NMR studies focus on adding a universal transformation on input waveform, we suspect that a more elaborate transformation is required for the music signal. Hence, we propose the following two adjustments.

First, music signals are fundamentally more complex than other audio data (e.g., speech command). To capture complicated harmonic relations, MIR tasks usually rely on perceptually meaningful features, such as Mel-spectrogram, instead of using raw audio. Moreover, the categorization of music is usually not determined by only one or two factors. In genre classification, for example, the same instrumentation or chord progression can appear in completely different music genres. Therefore, we propose that the trainable parameters should be related to features instead of the raw wave-

form. Moreover, the trainable parameters should be non-linear and “input-dependent” to capture the complex musical features of each training sample.¹ To achieve these criteria, we add a transformation function (T) composed of convolutional layers to transform the original feature, as shown in Fig 1b. The transformation function will learn how to add “noise” to the input features depending on each input sample. That is, $X' = \mathcal{H}(X; w) := T(X)$, where w is the parameters of T . We refer to this adjustment input-dependent NMR (*ID-NMR*).

Second, there are few potential drawbacks of directly modifying either waveform or input features. For example, to update parameters of the transformation function, the reprogramming layer still needs the gradient from the whole model during training, which leads to slower training time. Moreover, compared to the raw waveform or input features, middle layers of the pre-trained models sometime represent more high-level information which is more critical for classification. Hence, we propose using skip connections to transform the output of the middle layer v of the model, and add the transformed feature v' (or noise) to the classifier for classification, as shown in Figure 1c. By using skip connections, we expect the adaptive information from the input-dependent features from the skip-connection to directly influence the classifier. We will refer to this method later as *IDS-NMR*.

3.3. Source Models

In this section, we introduce two DNN models that we use in the experiment. To explore the probability of cross domain reprogramming, we choose SpeechATT and AST models pre-trained on speech and audio data, respectively.²

SpeechATT: we choose SpeechATT as the efficient baseline model since Yang et al. has achieved several state-of-the-art results on time-series data by using the NMR training scheme and the SpeechATT model. The model is pre-trained on the Google Speech Command dataset [5], which contains 105,829 utterances of 35 words. Different than other speech recognition tasks, utterances usually contain very short time frames, so the input of this model is only a one-second audio sample. In music genre classification tasks, the input audio is usually longer than one second. To satisfy the data format of pre-trained speech command model, we chunked the input sample into non-overlapping one-second pieces. The final probability of each label is attained by averaging the probabilities over all chunks.

Audio Spectrogram Transformer (AST): In recent years, transformer architectures have been successfully applied to a variety of tasks [4, 26]. In this paper, we experiment with the recently proposed AST model [18] for reprogramming. AST is a purely attention-based model for audio classification with BERT-alike patch-wise feature learning. To further

leverage a larger scale of training data, AST is pre-trained on ImageNet [3] and fine-tuned AudioSet [10]. Although other large-scale models, such as VGGish model [10], have been commonly used in several MIR downstream tasks [13, 12], we choose AST since it has outperformed the VGGish model on AudioSet classification. Moreover, reprogramming VGG-based architectures has shown worse performance due to the disadvantage of the multi-channel feature aggregation used in VGG-based architectures [16]. The input of AST model is an audio recording of length 10 s. We apply a similar chunking approach as for SpeechATT: audio is chunked into non-overlapping 10 s segments.

4. EXPERIMENT

4.1. Dataset

We choose GTZAN, a small but popular dataset for genre classification, for the experiment. It contains 10 musical genres, with each genre having 100 audio snippets of 30 s length, resulting in a total of 1000 snippets. Since the total hours of this dataset is only around 8 hours, it is suitable to represent the scenario of low-resource training data. We adopt the “fault-filtered” split from [27] which addresses some of the reported issues with this dataset [20], resulting in 443 pieces for training, 197 pieces for validation, and 290 pieces for testing. The same dataset split is also used in the MIREX baseline systems [12].

4.2. Experimental Setup

Four baseline methods are included in the experiment for comparison. The first baseline is a simple CNN architecture (*BL-CNN*). There are plenty of ways to build the baseline CNN. We choose ResNet since it has shown its efficiency on audio-tagging tasks [28]. We simply concatenate four ResNet blocks with kernel size 3 and channel size 50.³ For the second baseline, we simply fine-tune the original pre-trained model with the datasets mentioned above as another common TL method. We call this method *BL-FT-AST* in short. The above mentioned methods together with our approaches can fit the model on task-specific datasets.

Another common transfer learning method is to extract the representation from the last layer of the pre-trained model and train a shallow supervised model for classification. Although this method could not tune the pre-trained models on task-specific datasets, it is included for comparison. Following the work from [14], the mean pooling representations across time are used to train an one-layer MLPs with 512 hidden units for classification. Except for experimenting on the representation extracted from AST model (*BL-R-AST*), we also include the representation extracted from VGGish [10] model (*BL-R-VGGish*) for comparison. We use the Pytorch version of the VGGish model⁴ for this experiment.

¹We also experimented on transforming the waveform or adding noise to the features, but these methods resulted in low accuracy.

²The preprocessing step, DNN model and classifier all follow the implementation of the original models.

³The parameters are roughly tuned to produce the best result.

⁴<https://github.com/harritaylor/torchvggish>

Model	SpeechATT _{audio}	AST _{audio + vision}
II-NMR	0.399 ± 0.013	0.503 ± 0.018
ID-NMR	0.609 ± 0.005	0.802 ± 0.021
IDS-NMR	0.657 ± 0.021	0.828 ± 0.017

Table 1: Results on GTZAN dataset comparing different NMR methods and pre-trained models.

Method	Pre-trained Model	Accuracy
FT-CNN [7]	MusiCNN	0.790
FT-CNN [29]	SampleCNN	0.821
Emb.+SVM [12]	[7]+[30]	0.801
Probing [14]	JukeBox	0.797
BL-CNN	CNN _{audio}	0.666 ± 0.034
BL-R-VGGish	VGGish _{audio}	0.765 ± 0.017
BL-FT-AST	AST _{audio + vision}	0.772 ± 0.013
BL-R-AST	AST _{audio + vision}	0.831 ± 0.005
IDS-NMR	AST _{audio + vision}	0.828 ± 0.017

Table 2: Result on GTZAN dataset comparing existing works and other baseline methods with the *IDS-NMR* method.

Each setting is trained 100 epochs, and the model that performs the best on the validation set is picked for evaluation. Adam optimizer with 0.0001 learning rate is used for training. Label mapping layers have $n = 2$ for SpeechATT and $n = 5$ for AST after parameter search. Each method is trained five times with different random seeds. Following previous works [12, 29, 7, 14] we report the mean and standard deviation of accuracy as the evaluation metrics. For the convolutional layers (T), we use a similar architecture as the *BL-CNN*, with a kernel size of 3 and 136 channels, to roughly match the total parameters as the *BL-CNN* model. For the convolutional layers (T), we use a similar architecture as the *BL-CNN*, with a kernel size of 3 and 136 channels, to roughly match the total parameters as the *BL-CNN* model.

5. RESULT & DISCUSSION

We then compare the results between different NMR methods. The performance of different methods by reprogramming the SpeechATT and AST models can be seen in Table. 1. We can observe that for both models, directly adding trainable noise (*NMR-Noise*) does not work well. However, our proposed input-dependent methods can largely improve the performance. After adding skip-connection for training, both models achieve obvious performance gain. SpeechATT generally performs worse than AST since the original model is trained only one-second audio. For music genre classification, it is hard to capture context information within such a short audio clip. Moreover, SpeechATT has fewer parameters and is trained with a smaller and music-unrelated dataset.

In the end, we compare baseline methods and existing TL

Methods	# Trainable Params	Speed (min:sec)
BL-FT	88,132,063	-
II-NMR	160,000	10:08
ID-NMR	234,623	10:11
IDS-NMR	235,680	03:50

Table 3: Training speed and number of parameters by using different training methods on SpeechATT and AST model.

methods with our input-dependent skip connection method. Result can be seen in Table 2. We can observe that although our proposed method utilizes a model pre-trained on partially music-related datasets, AudioSet and ImageNet, it outperforms existing models pre-trained on music-specific datasets, such as MSD [7, 29], and model pre-trained on millions of music data [14]. Our *IDS-NMR* method also outperforms fine-tuning method (BL-FT-AST). Although our methods perform on par with the representation method (BL-R-AST), our *IDS-NMR* method gives models the flexibility to learn both mid-level and high-level information. As extra findings, our *IDS-NMR* method with jointly ImageNet and AudioSet pre-training could achieve one highest test accuracy (**85.1%**) and outperforms AudioSet-only pre-training by 2.8% within all the random seeds. The representation from AST model (*BL-R-AST*) performs better than VGGish model (*BL-R-VGGish*). This result suggests that the utilization of the AST model could thus lead to improved results for a variety of MIR applications in the near future, replacing VGGish features which have been commonly used in many music downstream tasks.

We also compare the number of parameters and training speed for the proposed methods. We use a single GeForce RTX 2080 GPU to train the GTZAN dataset and report how many seconds it takes to complete one epoch. In each epoch, the models visit each training sample once. We can see from Table 3 that although with or without skip connection has similar number of trainable parameters for our input-dependent methods, with skip connecting is faster in training. Fine-tuning AST model although can achieve good result, it could not fit on a single GPU for training.

6. CONCLUSION

In this work, we proposed a new TL training scheme to leverage the pre-trained models on downstream tasks. We prob into two pre-trained models, SpeechATT and AST, by using both baseline and improved NMR methods. Experiment result on low-resource music genre classification shows that the proposed input-dependent method not only outperforms the fine-tuning transfer learning method but also achieve higher accuracy than existing models pre-trained on music-specific datasets. In the future, we plan to explore the proposed methods on other MIR tasks which also suffer from data limitation problem. Our implementation will be available at <https://github.com/biboamy/music-repro>.

7. REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [2] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *Proc. CVPR*, 2009, pp. 248–255.
- [3] O. Russakovsky *et al.*, “Imagenet large scale visual recognition challenge,” *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [4] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in *Proc. NAACL-HLT*, 2019.
- [5] P. Warden, “Speech commands: A dataset for limited-vocabulary speech recognition,” *arXiv preprint arXiv:1804.03209*, 2018.
- [6] J. F. Gemmeke, D. P. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, “Audio set: An ontology and human-labeled dataset for audio events,” in *Proc. ICASSP*, 2017, pp. 776–780.
- [7] J. Pons and X. Serra, “musicnn: Pre-trained convolutional neural networks for music audio tagging,” in *Late-breaking/demo session in ISMIR*, 2019.
- [8] P. Dhariwal, H. Jun, C. Payne, J. W. Kim, A. Radford, and I. Sutskever, “Jukebox: A generative model for music,” *arXiv preprint arXiv:2005.00341*, 2020.
- [9] F. Korzeniowski, O. Nieto, M. McCallum, M. Won, S. Oramas, and E. Schmidt, “Mood classification using listening data,” in *Proc. ISMIR*, 2020, pp. 542–549.
- [10] S. Hershey *et al.*, “Cnn architectures for large-scale audio classification,” in *Proc. ICASSP*, 2017.
- [11] J. Cramer, H.-H. Wu, J. Salamon, and J. P. Bello, “Look, listen, and learn more: Design choices for deep audio embeddings,” in *Proc. ICASSP*, 2019, pp. 3852–3856.
- [12] K. C. B. Liang, Q. M. BU, and T. M. Entertainment, “Do user preference data benefit music genre classification tasks?” *MIREX*, 2020.
- [13] E. Koh and S. Dubnov, “Comparison and analysis of deep audio embeddings for music emotion recognition,” *AAAI Workshop on Affective Content Analysis*, 2021.
- [14] R. Castellon, C. Donahue, and P. Liang, “Codified audio language modeling learns useful representations for music information retrieval,” in *Proc. ISMIR*, 2021.
- [15] G. F. Elsayed, I. Goodfellow, and J. Sohl-Dickstein, “Adversarial reprogramming of neural networks,” in *Proc. ICLR*, 2018.
- [16] C.-H. H. Yang, Y.-Y. Tsai, and P.-Y. Chen, “Voice2series: Reprogramming acoustic models for time series classification,” in *Proc. ICML*, 2021, pp. 11 808–11 819.
- [17] R. Vinod, P.-Y. Chen, and P. Das, “Reprogramming language models for molecular representation learning,” *arXiv preprint arXiv:2012.03460*, 2020.
- [18] Y. Gong, Y.-A. Chung, and J. Glass, “AST: Audio spectrogram transformer,” in *Proc. Interspeech*, 2021.
- [19] D. C. de Andrade, S. Leo, M. L. D. S. Viana, and C. Bernkopf, “A neural attention model for speech command recognition,” *arXiv preprint arXiv:1808.08929*, 2018.
- [20] B. L. Sturm, “The GTZAN dataset: Its contents, its faults, their effects on evaluation, and its future use,” *arXiv preprint arXiv:1306.1461*, 2013.
- [21] I. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” in *Proc. ICLR*, 2015.
- [22] Y.-Y. Tsai, P.-Y. Chen, and T.-Y. Ho, “Transfer learning without knowing: Reprogramming black-box machine learning models with scarce data and limited resources,” in *Proc. ICML*, 2020, pp. 9614–9624.
- [23] K. Hambarzumyan, H. Khachatrian, and J. May, “Warp: Word-level adversarial reprogramming,” in *Proc. ACL-IJCNLP*, 2021, pp. 4921–4933.
- [24] K.-W. Chang, W.-C. Tseng, S.-W. Li, and H.-y. Lee, “An exploration of prompt tuning on generative spoken language model for speech processing tasks,” *arXiv*, 2022.
- [25] H. Yen, P.-J. Ku, C.-H. H. Yang, H. Hu, S. M. Siniscalchi, P.-Y. Chen, and Y. Tsao, “A study of low-resource speech commands recognition based on adversarial reprogramming,” *arXiv preprint arXiv:2110.03894*, 2021.
- [26] N. Parmar, A. Vaswani, J. Uszkoreit, L. Kaiser, N. Shazeer, A. Ku, and D. Tran, “Image transformer,” in *Proc. ICML*, 2018, pp. 4055–4064.
- [27] C. Kereliuk, B. L. Sturm, and J. Larsen, “Deep learning and music adversaries,” *IEEE Trans. on Multimedia*, vol. 17, no. 11, pp. 2059–2071, 2015.
- [28] M. Won, A. Ferraro, D. Bogdanov, and X. Serra, “Evaluation of cnn-based automatic music tagging models,” in *Proc. SMC*, 2020.
- [29] J. Lee, J. Park, K. L. Kim, and J. Nam, “Samplecnn: End-to-end deep convolutional neural networks using very small filters for music classification,” *Applied Sciences*, vol. 8, no. 1, p. 150, 2018.
- [30] P. Covington, J. Adams, and E. Sargin, “Deep neural networks for youtube recommendations,” in *Proc. ACM RecSys*, 2016, pp. 191–198.