# TOWARDS EXPLAINING EXPRESSIVE QUALITIES IN PIANO RECORDINGS: TRANSFER OF EXPLANATORY FEATURES VIA ACOUSTIC DOMAIN ADAPTATION

*Shreyan Chowdhury*[1]    *Gerhard Widmer*[1,2]

[1]Institute of Computational Perception and [2] LIT AI Lab, Johannes Kepler University Linz, Austria

## ABSTRACT

Emotion and expressivity in music have been topics of considerable interest in the field of music information retrieval. In recent years, *mid-level perceptual features* have been suggested as means to explain computational predictions of musical emotion. We find that the diversity of musical styles and genres in the available dataset for learning these features is not sufficient for models to generalise well to specialised acoustic domains such as solo piano music. In this work, we show that by utilising unsupervised domain adaptation together with receptive-field regularised deep neural networks, it is possible to significantly improve generalisation to this domain. Additionally, we demonstrate that our domain-adapted models can better predict and explain expressive qualities in classical piano performances, as perceived and described by human listeners.

***Index Terms***— Music, expressivity, domain adaptation

## 1. INTRODUCTION

*Domain mismatch* – a discrepancy between the kind of data available for training a classifier and the data on which it should then operate – is an important real-world problem, also in the field of acoustic recognition. For instance, the DCASE 2019 and 2020 challenges had dedicated tasks on *Acoustic Scene Classification with multiple/mismatched recording devices*.[1] The machine learning answer to this problem is research on effective methods for *transfer learning* and (supervised and unsupervised) *domain adaptation*.

The work presented in this paper is motivated by a particularly difficult acoustic transfer problem involving a complex musical phenomenon. In a large project,[2] we aim at studying the elusive concept of *expressivity in music* with computational and, specifically, machine learning methods. One aspect of that is the art of *expressive performance*, the subtle, continuous shaping of musical parameters such as tempo, timing, dynamics, and articulation by experienced musicians, while playing a piece, in this way imbuing the piece with

particular expressive and emotional qualities [1]. The *Con Espressione Game* was a large-scale data collection effort we set up in order to obtain personal descriptions of perceived *expressive qualities*, with the goal of studying human perception and characterisation of expressive aspects in performances of the same pieces by different artists [2].

In analyzing this data, we are now interested in seeing whether these subjective characterisations of expressive qualities are consistent and systematic enough for a machine to be able to predict them – at least partially – from the audio recordings. Moreover, we aim at obtaining musical insights: we want *interpretable* models that point to specific musical qualities that might underlie perceived expressive qualities. A set of musical descriptors that seem particularly suited to this was proposed in [3], where *mid-level musical features* were described that are intuitively understandable to the average musical listener, and a corresponding human-annotated set of music recordings was published (see next section). In [4], we had shown how such mid-level features, predicted from audio via trained classifiers, can be exploited to provide intuitive explanations in the context of emotion and mood recognition in general (non-classical) music. This was extended to a two-level explanation scheme in [5] which permitted to trace the mid-level explanations back to properties of the acoustic signal. There is reason to believe that some of these features may also hold predictive and explanatory power for expressive aspects in piano performance.

This is where a severe *mismatch problem* arises: there is no annotated ground truth data available for training mid-level feature extractors in classical piano music, and obtaining such data would be extremely cumbersome. At the same time, recordings of solo piano music are very different, musically and acoustically, from the kind of rock and pop music contained in the available mid-level training dataset. It is thus likely that a classifier trained on the latter will not generalise well to our piano recordings.[3]

In this paper, we present several steps to bridge this domain mismatch through architecture choice and unsupervised domain adaptation techniques, and show that they are effec-

---

[1]e.g., http://dcase.community/challenge2019/task-acoustic-scene-classification-results-b

[2]https://www.jku.at/en/institute-of-computational-perception/research/projects/con-espressione

---

[3]Note that we cannot test this directly, as we have no mid-level feature ground truth for the *Con Espressione* performances. (We will use the few piano recordings in the midlevel dataset as our domain adapatation test set – see Section 4.)

tive in generalising a model to solo piano recordings. In a final step, we will try the adapted classifier on the Con Espressione recordings, testing whether domain adaptation improves the predictability of expressive qualities from mid-level features predicted from audio, and identifying those features that seem to have specific predictive and explanatory power.

## 2. MID-LEVEL FEATURES AND THE CON ESPRESSIONE DATA

### 2.1. The Mid-level Features Dataset

Seven mid-level musical features were proposed in [3], viz. *melodiousness, articulation, rhythmic complexity, rhythmic stability, dissonance, tonal stability*, and *modality* (or "minorness"). To approximate these perceptual features for a set of audio clips, the authors took a data-driven approach. Ratings from listeners were modelled into the final feature values that were made available as labels in the associated dataset (which we call the *Mid-level Features Dataset*) along with the audio clips. The labels for each feature are in the form of continuous values between 1 and 10 (the learning task for our models will thus be a regression task.) The exact questions asked to the listeners for rating each perceptual feature can be found in [3]. The audio clips chosen for the dataset come from different sources such as `jamendo.com`, `magnatune.com`, and the Soundtracks dataset [6]. There are a total of 5,000 clips of 15 seconds each in the dataset.

### 2.2. The *Con Espressione Game* Dataset

In the *Con Espressione Game*, participants listened to extracts from recordings of selected solo piano pieces (by composers such as Bach, Mozart, Beethoven, Schumann, Liszt, Brahms) by a variety of different famous pianists (for details, see [2]) and were asked to describe, in free-text format, the *expressive character* of each performance. Typical characterisations that came up were adjectives like "cold", "playful", "dynamic", "passionate", "gentle", "romantic", "mechanical", "delicate", etc. From these textual descriptors, the authors obtained, by statistical analysis of the occurrence matrix of the descriptors, four underlying continuous expressive dimensions along which the performances can be placed. These are the (numeric) target dimensions that we wish to predict via the route of mid-level features predicted from the audio recordings.

The *central challenge* in this is that the Mid-level Features Dataset [3], consisting mainly of pop, rock, hip-hop, jazz, electronic, and film soundtrack music, is vastly different, in sound and musical style, from the music of the Con Espressione dataset. This results in what is known as a *covariate shift* [7] between the training and the testing data.

In the following section, we describe a deliberate choice of training architecture that results in better generalisability of the trained models, and then present a two-step method to further adapt the model to our domain of choice.

## 3. MID-LEVEL FEATURE LEARNING VIA DOMAIN ADAPTATION (DA)

In the following sections, *target domain* refers to solo piano performance audio and *source domain* refers to all other musical audio (non-piano audio clips in the Mid-level Features Dataset).

### 3.1. Receptive Field Regularized ResNet

As a first step towards improving out-of-domain generalization of mid-level feature prediction, we switch from the VGG-ish network of [4] to the Receptive-Field Regularized ResNets (RF-ResNet) originally introduced in [8] for acoustic scene classification and later shown to work well for music information retrieval tasks as well [9]. The rationale behind this is that the smaller receptive field of the RF-ResNet prevents overfitting, particularly when the training data is limited in quantity. The architecture differs from a regular ResNet [10] by reducing the kernel sizes of several convolutional layers and removing some max pooling layers. Our RF-ResNet consists of three stages with three residual blocks in the first stage and one residual block each in the second and third stages. The last stage consists of only 1-by-1 convolutional layers. There are two max pooling layers in the first stage between the convolutional blocks, and one average pooling layer after the third stage before going into a final 1-by-1 convolutional feed forward layer. The output is a seven-dimensional vector where the elements correspond to the predictions of each of the seven mid-level features.

### 3.2. Unsupervised DA through Backpropagation

We adopt the *reverse-gradient* method introduced in [11], which achieves domain invariance by adversarially training a domain discriminator attached to the network being adapted, using a gradient reversal layer. The procedure requires a large unlabelled dataset of the target domain in addition to the labelled source data. The discriminator tries to learn discriminative features of the two domains but due to the gradient reversal layer between it and the feature extracting part of the network, the model learns to extract domain-invariant features from the inputs.

This adaptation procedure is applied to the RF-ResNet described above. Since our target domain of interest solo piano performance music, we use audio from the MAESTRO dataset [12] as our unlabelled data source. It contains more than 200 hours of recorded piano performances. During training, each batch that the model sees contains an equal number of labelled source data points and unlabelled target data points. The regressor/classifier head of the model tries to predict the source labels while the discriminator head predicts the domain for each data point in the batch. The combined loss of the two heads is then backpropagated while reversing the gradient after the discriminator during the backward pass.

### 3.3. Teacher-Student Training Scheme

As a final step, we refine our domain adaptation using a teacher-student training scheme tailored to our scenario (see Fig.1). We train multiple domain-adaptive models using the unsupervised DA method of Section 3.2 and use these as teacher models that are eventually used to assign pseudo-labels to our unlabelled MAESTRO dataset. Before the pseudo-labelling step, we select the best performing teacher models with the validation set. Even though the validation set contains data from the source domain, this step ensures that models with relatively lower variance are used as teachers. This helps filter out the particularly poorly adapted models from the previous step, which may occur due to the inherently less stable nature of adversarial training methods [13].

After selecting a number of teacher models (in our experiments, we used four), we label a randomly selected subset of our unlabelled dataset using predictions aggregated by taking the average. This pseudo-labelled dataset is combined with the original labelled source dataset to train the student model. We observed that the performance on the test set increased until the pseudo-labelled dataset was about 10% of the labelled source dataset in size, after which it saturated.

The teacher-student scheme allows the collective "knowledge" of an ensemble of adapted networks to be distilled into a single student network. The idea of knowledge distillation, which was originally introduced for model compression in [14], has been used for domain adaptation in a supervised setting previously in [15]. The distillation process functions as a regularizer resulting in a student model with better generalisability than any of the individual teacher models alone. Additionally, it can be thought of as a stabilisation step helping to filter out the adversarially adapted models that result from non-optimal convergence.

## 4. EXPERIMENTAL RESULTS

Since we have no ground truth labels for our real data of interest (the classical piano music) to evaluate the domain adaptation experiments[4], we created a ("piano"/target) test set manually by selecting clips from the Mid-level Features Dataset containing only solo piano. This resulted in a set of 79 piano clips from the total of 5000. The other 4921 clips ("non-piano"/source) were split into training (90%), validation (2%) and test (8%) sets such that the artists in these sets are mutually exclusive (following [3]). The validation set is used to tune hyperparameters and for early stopping.

The inputs to all our models were log-filtered spectrograms (149 bands) of 15-second audio clips sampled at 22.05 kHz with a window size of 2048 samples and a hop length of 704 samples, resulting in $149 \times 469$-sized tensors. For training, we use the Adam optimizer with a multi-step learning rate scheduler. In the unsupervised DA step, the recordings
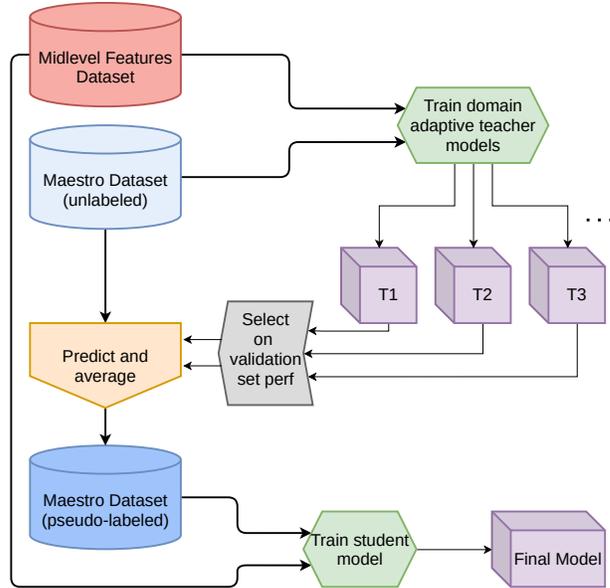
**Fig. 1**. Teacher-Student training scheme for unsupervised domain adaptation.

from the MAESTRO dataset are split into 15-second segments and a random subset of size equal to the mid-level training set is sampled on each run. During the pseudo-labelling stage, a random subset of 500 segments is sampled.

We observe (Fig. 2) that each of the steps mentioned in the previous section results in an improvement in the performance on the "piano" test set without compromising the performance on the "non-piano" one. In fact, we see a slight improvement in the non-piano metric upon introducing DA. This could be due to the presence of some data points similar to the target domain – for instance excerpts from piano concertos, which are not included in the "piano" test set.

To investigate our results further, we look at the discrepancy between the source and target domains in the representa-
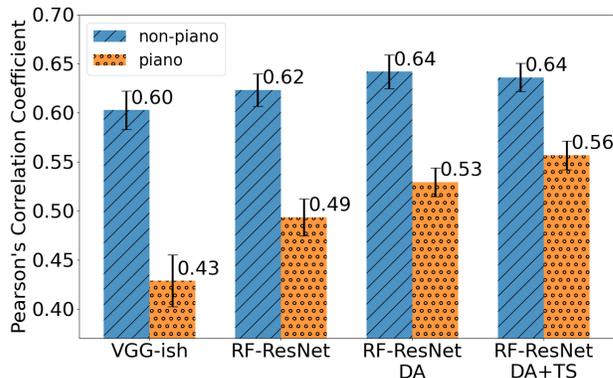


**Fig. 2**. Performance of mid-level feature models on non-piano and piano test sets.
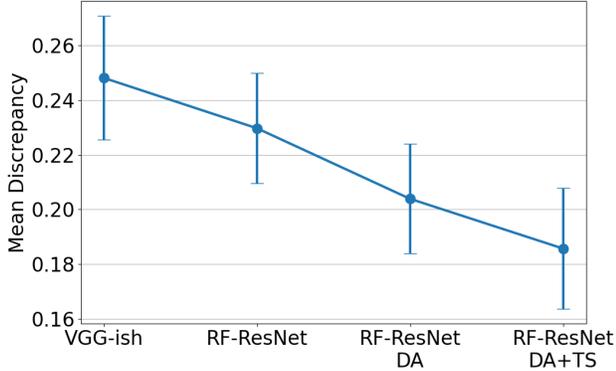
**Fig. 3**. Mean discrepancy between piano and non-piano sets.

|  | Dim 1 | Dim 2 | Dim 3 | Dim 4 |
|---|---|---|---|---|
| VGG-ish | 0.35 | 0.10 | 0.22 | 0.32 |
| RF-ResNet | 0.36 | 0.07 | 0.28 | 0.33 |
| RF-ResNet DA | **0.40** | 0.09 | **0.29** | 0.32 |
| RF-ResNet DA+TS | 0.35 | **0.15** | **0.29** | **0.34** |

**Table 1**. Coefficient of determination ($R^2$-score) of description embedding dimensions of the Con Espressione game using a linear regressor trained on predicted mid-level features.

| RF-ResNet | | RF-ResNet DA+TS | |
|---|---|---|---|
| Feature | $r$ | Feature | $r$ |
| articulation | 0.47 | melodiousness | $-0.39$ |
| rhythmic complexity | 0.41 | articulation | 0.46 |
| | | rhythmic complexity | 0.41 |
| | | dissonance | 0.40 |

**Table 2**. Pearson's correlation ($r$) for mid-level features with the first description embedding dimension, with (right) and without (left) domain adaptation. Features with $p < 0.05$ and $|r| > 0.20$ are selected. This dimension has positive loadings for words like "hectic", "irregular", and negative loadings for words like "sad", "gentle", "tender".

tion space, since it is known that the performance of a model on the target domain is bounded by this discrepancy [7]. We use the method given in [16] to compute the empirical distributional discrepancy between domains for a trained model $\phi$, which is given as $D(S', T'; \phi)$ in Eq. 1:

$$D(S', T'; \phi) = \left\| \frac{1}{m} \sum_{x \in S'} \phi(x) - \frac{1}{n} \sum_{x \in T'} \phi(x) \right\|_2 \quad (1)$$

where $S'$ is a population sample of size $m$ from the source domain and $T'$ is a population sample of size $n$ from the target domain. We observe that the discrepancy decreases for each step (Fig.3), justifying our three-step approach and explaining the improvement in performance.

## 5. PUTTING IT TO THE TEST

As a final step, we now return to our real target domain of interest and briefly investigate whether our domain-adapted models can indeed predict better mid-level features for modelling the expressive descriptor embeddings of the Con Espressione dataset. We do this by predicting the average mid-level features (over time) for each performance using our models and training a simple linear regression model on these features to fit the four embedding dimensions. Even though this is a very abstract task, for a variety of reasons – the noisy and varied nature of the human descriptions; the weak nature of the numeric dimensions gained from these; the complex and subjective nature of expressive music performance – it can be seen that the features predicted using domain-adapted models give comparatively better $R^2$-scores for all four dimensions.

Taking a closer look at Dimension 1 – the one that came out most clearly in the statistical analysis of the user responses and was characterized by descriptions like "hectic" and "agitated" (as opposed to, e.g., "calm" and "tender"; see [2]) – and looking at the individual mid-level features (see Table 2), we find that, first of all, the predicted features that show a strong

correlation with this dimension do indeed make sense: one would expect articulated ways of playing (e.g., with strong *staccato*) and rhythmically complex or uneven playing to be associated with an impression of musical agitation. What is more, after domain adaptation, the set of explanatory features grows, now also including perceived dissonance as a positive, and perceived melodiousness of playing as a negative factor – which again makes musical sense and testifies to the potential of domain adaptation for transferring explanatory acoustic and musical features.

## 6. CONCLUSION

In this paper, we presented a three-step approach to adapt mid-level models for recordings of solo piano performances. We significantly improved the performance of these models on piano audio by using a receptive field regularised network and performing unsupervised domain adaptation via a teacher-student training scheme. We also demonstrated improved prediction of meaningful perceptual features corresponding to expressive dimensions. We conclude that this route of domain adaptation shows potential for a more general task of adapting models to specific genres or musical styles.

## 7. ACKNOWLEDGMENT

# 8. REFERENCES

[1] Carlos E Cancino-Chacón, Maarten Grachten, Werner Goebl, and Gerhard Widmer, "Computational Models of Expressive Music Performance: A Comprehensive and Critical Review," *Frontiers in Digital Humanities*, vol. 5, pp. 25, 2018.

[2] Carlos Cancino-Chacón, Silvan Peter, Shreyan Chowdhury, Anna Aljanaki, and Gerhard Widmer, "On the Characterization of Expressive Performance in Classical Music: First Results of the Con Espressione Game," in *Proceedings of the 21st International Society for Music Information Retrieval Conference (ISMIR)*, 2020.

[3] Anna Aljanaki and Mohammad Soleymani, "A Data-driven Approach to Mid-level Perceptual Musical Feature Modeling," in *Proceedings of the 19th International Society for Music Information Retrieval Conference, (ISMIR)*, 2018, pp. 615–621.

[4] Shreyan Chowdhury, Andreu Vall, Verena Haunschmid, and Gerhard Widmer, "Towards Explainable Music Emotion Recognition: The Route via Mid-level Features," in *Proceedings of the 20th International Society for Music Information Retrieval Conference (ISMIR)*, Delft, The Netherlands, 2019.

[5] Verena Haunschmid, Shreyan Chowdhury, and Gerhard Widmer, "Two-level Explanations in Music Emotion Recognition," in *International Conference on Machine Learning (ICML)*, 2019, Machine Learning for Music Discovery workshop.

[6] Tuomas Eerola and Jonna K. Vuoskoski, "A Comparison of the Discrete and Dimensional Models of Emotion in Music," *Psychology of Music*, vol. 39, no. 1, pp. 18–49, 2011.

[7] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan, "A Theory of Learning from Different Domains," *Machine learning*, vol. 79, no. 1-2, pp. 151–175, 2010.

[8] Khaled Koutini, Hamid Eghbal-Zadeh, Matthias Dorfer, and Gerhard Widmer, "The Receptive Field as a Regularizer in Deep Convolutional Neural Networks for Acoustic Scene Classification," in *2019 27th European signal processing conference (EUSIPCO)*. IEEE, 2019, pp. 1–5.

[9] Khaled Koutini, Shreyan Chowdhury, Verena Haunschmid, Hamid Eghbal-zadeh, and Gerhard Widmer, "Emotion and Theme Recognition in Music with Frequency-Aware RF-Regularized CNNs," in *Multimedia Evaluation Benchmark (MediaEval) 2019 Workshop*, 2019.

[10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep Residual Learning for Image Recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[11] Yaroslav Ganin and Victor Lempitsky, "Unsupervised Domain Adaptation by Backpropagation," in *International Conference on Machine Learning (ICML)*. PMLR, 2015, pp. 1180–1189.

[12] Curtis Hawthorne, Andriy Stasyuk, Adam Roberts, Ian Simon, Cheng-Zhi Anna Huang, Sander Dieleman, Erich Elsen, Jesse Engel, and Douglas Eck, "Enabling Factorized Piano Music Modeling and Generation with the MAESTRO Dataset," in *International Conference on Learning Representations*, 2019.

[13] Tong Che, Yanran Li, Athul Paul Jacob, Yoshua Bengio, and Wenjie Li, "Mode Regularized Generative Adversarial Networks," *arXiv preprint arXiv:1612.02136*, 2016.

[14] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean, "Distilling the Knowledge in a Neural Network," *arXiv preprint arXiv:1503.02531*, 2015.

[15] Taichi Asami, Ryo Masumura, Yoshikazu Yamaguchi, Hirokazu Masataki, and Yushi Aono, "Domain Adaptation of DNN Acoustic Models using Knowledge Distillation," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 5185–5189.

[16] Yu Sun, Eric Tzeng, Trevor Darrell, and Alexei A Efros, "Unsupervised Domain Adaptation through Self-Supervision," *arXiv preprint arXiv:1909.11825*, 2019.