

Figure 2. Proposed system for multimodal classification from pose and audio time series extracted from 12s video examples. This diagram shows 5 such configurations. A and B represent unimodal classification for video and audio streams respectively. C and D represent multimodal classification using data fusion at different layers. E represents model fusion using unimodal classification results.

[17, 18]. An important analysis parameter (that also dictates the window size) is the pitch search range, limiting which can help minimize octave errors. The tonic, automatically detected and manually verified [19, 20], was obtained for each performance and used to achieve the needed tonic-normalization of the extracted pitch contours. The detected tonic is also used to define the expected 2 octave pitch range for each piece. The voicing is a binary variable per frame that is set to zero in detected silence frames corresponding to singing pauses. Brief unvoiced regions (less than 500 ms) arising from short silences and consonant utterances are filled in via cubic spline interpolation to obtain the continuous pitch contours associated with melodic movements. The resulting time series from each clip thus comprises of 1200 samples (12s x 100/s).

3.2 Video Features

Real-time skeleton data, such as that obtained via OpenPose, can be handled by a graph model such as Graph CNN. We treat it instead as multivariate time-series data, each time series corresponding to one of the position coordinates of a tracked joint, similar to the processing of data from body-worn sensors in human activity recognition tasks. The 2D positions of the 11 keypoints of the upper body (eyes, nose, neck, shoulders, mid-hip, elbows and wrists) are recorded and used to obtain normalised (x,y) coordinates for each of the two wrists at the frame rate of 25/s. The position data for the singer’s two wrists alone is retained for the analysis since (a) these points are more reliably estimated by OpenPose than others such as the elbows or shoulders, and (b) hand gestures most clearly relate to raga expression, and this data is most likely to contribute to raga classification. We thus have 4 time-series representing the x,y positions of each wrist. Any missing data are interpolated and each of the wrist position time series is low-pass filtered to remove any jitter. The length of each of the time series is 300 samples (12s x 25 fps).

3.3 Network Architectures and Hyperparameters

Given that our music audio and video time series data embed information at multiple time scales, we choose an inception network for its multiple kernel sized filters [21,22]. Inception networks have been previously used for multivariate time series classification [23, 24]. We empirically observed that preceding the inception block with two convolutional layers led to superior audio performances while a single convolutional layer worked best for the video input. The relatively high frame rate of audio also necessitated greater stride choices through the convolutional layers that helped reduce the time series dimension to 200 at the input to the inception block. The convolutional layers help in learning audio features which may be relevant for processing via the inception network which further has a range of receptive fields for the convolution. A similar approach to use prior convolution layers with inception blocks was adopted in the original work introducing the inception network for image classification [21]. We further exploit these prior convolutional layers learned for the individual unimodal (audio and video) channels in our latent representation fusion model described in the following section. The inception block is followed by pooling and softmax layers.

Figure 3 shows the architecture of the inception block. Overall, the hyperparameters tuned include the kernel size and number of filters in the convolution layers, the common kernel size, number of filters, pool size and pool type in the inception block and the pool type in the final pooling layer. Hyperparameter tuning was carried out with sweeping across the chosen ranges using Bayesian optimization [25, 26].

3.4 Multimodal Analysis

It is widely believed that integrating features from multiple modalities can lead to more robust classification due

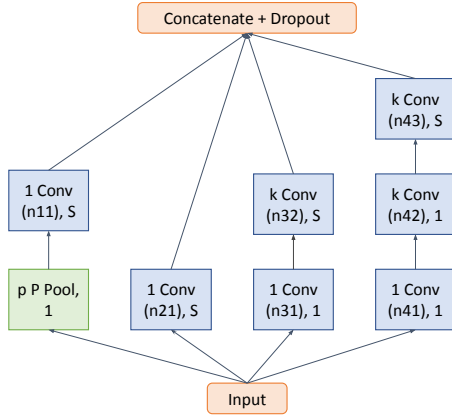


Figure 3. Structure of the inception block. ‘k Conv (n), S’ indicates a convolution layer with n k-sized kernels and a stride S. ‘p’ indicates the pooling size of the pool layer and ‘P’ the type of pooling. k, p, P and the number of filters were determined by hyperparameter tuning. S=1 for the video and combined models and S=2 for the audio model.

to potentially complementary information across the different modalities, all observing the same phenomenon. Multimodal classification has been an active research area in machine learning [27] where the precise approaches to combining information have been broadly categorised as early (feature-based), late (decision-based) and hybrid (their combination). While late fusion builds on combining the decisions of the individual unimodal predictors, early fusion can potentially exploit the low-level correlations between features across modalities. In our task, we have audio and video time series that actually co-vary as the singer makes continuous movements with her hands while singing. This relationship between melodic pitch variation and wrist movements makes it attractive to investigate the combination of the time series for classification.

In Figure 2, the flow-graphs labeled A and B depict the individual video and audio classification paths respectively. While the basic audio-only and video-only models differ only in the number of convolutional layers employed before the inception block, the hyperparameters assume values that are influenced by the time resolution of the corresponding time series with their different sampling rates. Early fusion is then obtained by first downsampling the pitch and voicing time series to the lower rate of 50/s and interpolating the wrist position data from its original 25/s to 50/s. The six time series form a 6-channel input to the convolutional layers thus realizing a form of source fusion. This is shown in position C in Figure 2 and a separate convolutional layer is used to learn the joint features before passing them onto the inception block.

Late fusion is achieved at the point labeled E (decision stage) in Figure 2 by combining the softmax outputs of the ensemble of the best model for each modality. We use soft voting [28] between the two classifiers by averaging the softmax outputs and then choosing the class with the maximum average as the predicted class. In addition, we learn

Data split	Seen split			Unseen split		
	AG	CC	SCh	AG	CC	SCh
Train	5590	5487	5588	4715	4105	4304
Validation	972	1075	974	1847	2457	2258

Table 3. Number of 12 s duration examples in each singer’s train-validation data split

classifiers using the concatenated softmax outputs of the best models of each modality. This is a common approach in model stacking [29,30] and has been used in multimodal fusion earlier [27,31]. We try multiple models for stacking viz. Random Forests [32,33] and logistic regression [34] and hyperparameter tune each of them using scikit-learn’s GridSearchCV routine [35].

While decision fusion is the most straightforward approach to combining modalities, we also consider fusion at an intermediate stage given that the dynamic correspondence between the movement in the video and the aligned audio is expected to persist in the earlier convolutional layers. We achieve this by fusing the latent representations from each of the subnetworks where each has been optimised for a simpler task, viz. audio-only or video-only classification. Here we use the pre-trained weights of the convolutional layers of the best models of each modality and freeze the weights. This ensures that the features learnt for the individual modalities are maintained and the following inception block is tuned to amplify the inter-relations between the two modalities. We call this novel method of fusion as latent fusion and depict it by D in Figure 2. Given that the input audio and video streams are at different effective temporal rates (lengths of 200 and 300 respectively) due to their original sampling and/or the different convolutional layer strides, we need to do a pooling on each individual modality convolutional output to bring the number of steps in the output sequence of the frozen models in sync. In addition, post the fusion we apply a convolution with filter size of 1 and learn the number of required filters via hyperparameter tuning thus letting the model learn the number of channels needed for the inception network.

4. EXPERIMENTS AND RESULTS

The training and validation sets are designed separately for the two tasks. We report experimental results for different model architectures (including one or both modalities) for each singer. In the seen singer task, we create training and validation data splits for each singer. Raga classification accuracies are reported for each method on the validation dataset of each singer and also averaged across the singers. For a given singer, the validation set comprises the set of examples from one of the singer’s alap takes for each raga. The corresponding training dataset is then all the *remaining* pieces by the singer plus all the pieces of the other 2 singers. We have thus a validation and train set for each of the 3 singers. The similar exercise is carried out for the unseen singer task. Here, all the pieces by a given singer

are placed in the validation set and the pieces sung by the other two singers in the training set. Table 3 summarises the number of 12 s examples in each split as used in the experiments presented next.

For each task, we carry out model hyperparameter tuning on each individual singer’s train-val dataset to obtain 3 sets of hyperparameters in all. These 3 hyperparameter-tuned models next have their weights recomputed on the training data of a given singer (say, AG) to obtain 3 trained models. Finally, the evaluation of Singer AG validation data is obtained via the ensemble of the three models. Similarly, we obtain the ensembled models for each of the singers CC and SCh and evaluate the methods on their respective validation datasets.

Our first set of experiments tests separately the performances of the audio modality and the video modality for the seen and unseen singer raga classification tasks. Table 4 presents the obtained validation accuracies. We observe that the unseen singer task is more challenging as expected. The audio based accuracies are significantly higher than those from video data on both tasks. The video based accuracies in the unseen singer task are close to chance (in the 9-way raga classification) indicating that the association between raga identity and gesture is highly singer dependent.

Our next set of experiments involve different approaches to multimodal classification. Given the non-informative nature of video cues in the unseen singer task, we restrict our attention here to the seen singer task. Table 5 presents the obtained validation accuracies across the different classification methods for each singer and the resulting mean across singers. We observe that early fusion is on the average at the performance level of the weaker modality. This is similar to the observations of Oramas [36] where the learning over very unequal modalities can be overwhelmed by either one. The results of late fusion appear in the final two rows and we find that they fall slightly short of the audio performances, all pointing to the challenge of actually realizing the benefits of the complementary information.

We look for another opportunity to combine audio and video information streams at the output of the convolutional layers. The latent representations at this stage are generated from convolutional layers that are frozen in pre-trained unimodal classification tasks. Rather than simple concatenation of the two representations, we use pooling to first align the representations along the time axis reducing both audio and video to a length of 100 from 200 and 300 respectively. We obtain a performance that is slightly, but consistently, superior to that from audio alone (D vs B in Table 5) for every one of the 3 singers, indicating that the combination of aligned latent representations successfully exploits the joint information. This is also borne out by the histogram summary provided in Figure 4.

5. DISCUSSION

As expected, prediction accuracy is significantly higher for the audio modality than the video, and higher in the ‘seen

Data Split	Seen Singer		Unseen Singer	
	Audio	Video	Audio	Video
AG	92.1	36.3	76.9	14.3
CC	79.4	31.8	60.4	13.8
SCh	77.0	39.2	67.2	10.0

Table 4. Validation accuracy (%) of only audio and only video modalities on seen/unseen singer train-val splits.

Model Type	Model Name	AG	CC	SCh	Mean
A	Video	36.3	31.8	39.2	35.8
B	Audio	92.1	79.4	77.0	82.8
C	Source fusion	30.1	42.4	35.8	36.1
D	Latent fusion	93.3	82.7	79.2	85.1
E1	Equal voting	85.9	73.7	67.9	75.8
E2	Stacking classifier – RF	81.9	74.2	76.3	77.5

Table 5. Validation accuracy (%) from each singer’s split for different model architectures in the seen singer task.

singer’ than the ‘unseen singer’ condition. Prediction accuracy is increased slightly when audio and video data are combined in comparison with audio data alone (from 82.8 to 85.1 %). It should be noted that some extracts score much more highly than others: SCh’s Shree scores highly for prediction accuracy in both modalities; CC’s MM is unusual in that video prediction seems more accurate than the audio; some extracts score very poorly on prediction from video only, such as AG’s Kedar and SCh’s Bahar. Some of the wrong predictions can probably be attributed to lack of data: for example, if the singer is silent for a significant part of the 12 second clip, or their hands do not move significantly. Other mismatches may be explained by features of either the melodic movement or the singers’ gestures.

We present confusion matrices for the audio, video and the latent fusion bimodal classification in Figure 5. The video-only matrix has significant off-diagonal dispersion.

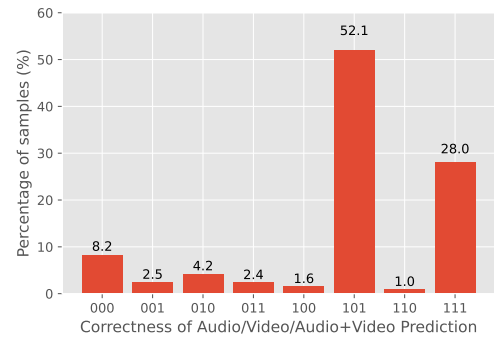


Figure 4. Histogram indicating percentage of the 3021 validation data examples predicted correctly (1) or incorrectly (0) by audio, video and the latent fusion based methods. For example, 011 indicates incorrect prediction by audio but correct predictions by video and bimodal classifiers.

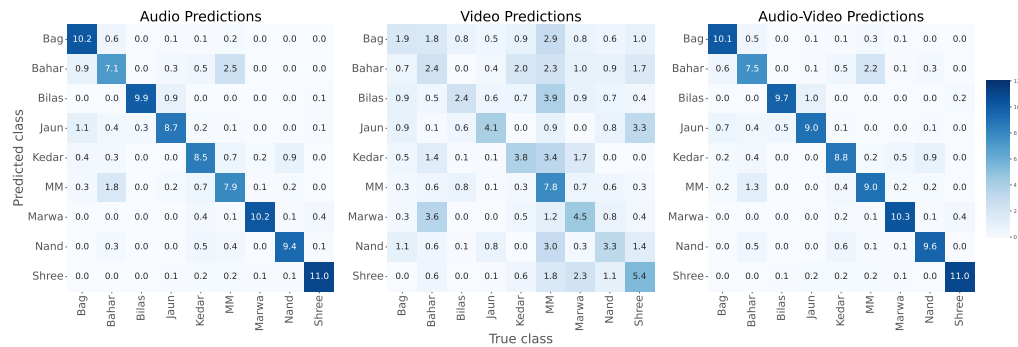


Figure 5. Confusion matrices of predictions made from audio, video and audio-video modalities. Numbers are represented in percentages of the total number of test examples across the three singers combined.

On the other hand, the multimodal matrix visibly improves upon the audio-only by further moving examples into the diagonal. The similar behaviour was noted in the individual singer confusion matrices. We also observe that the most common mis-predictions in the audio modality occur when the scale is the same or similar, and particularly where both pitch material and melodic movements are similar, as is the case between Miyan ki Malhar (MM) and Bahar. The main differences between these two closely-related ragas is that Bahar favours a higher pitch range and faster movement [37]: qualitative review of the prediction results shows that these factors are at play (e.g., in the faster portions of the extracts MM is more likely to be mis-classified as Bahar, at least for the two female singers). Other confusions are observed between Kedar and Nand (same scale), between Bageshree and Bahar (the latter uses one additional note), or between Jaunpuri and Bageshree (one note is different). Relatively few cases occur where the pitch material is very different (e.g. in Sch between MM and Nand or Kedar), which are harder to interpret. Some prediction errors seem to be related to the use of a low pitch (below Pa, the fifth degree, in the lower octave), as for example CC’s Shree between 38-53s: this can be explained by the fact that the pitch extraction was constrained to a range of two octaves, with the low Pa as its limit. We find in this case clear instances of video and multimodal predictions doing well.

Confusion matrices for the video prediction are harder to read, given the lower accuracy overall, but some patterns can be observed. The Bahar/MM confusion is also present in the video domain (for the female singers), perhaps related to the fact that these two ragas share not only a scale but also many aspects of their melodic movement. Bahar can sometimes be confused with Nand and Marwa. The Bahar/Nand confusion is not surprising, since both are associated with lively, complex melodic movement and a joyful mood. The confusion between these two ragas and Marwa is less expected, since Marwa’s mood contrasts strongly (being serious and often described as unsettled or restless). It may be that there is a similarity in the singers’

hand movements between Bahar and Nand’s liveliness and Marwa’s restlessness. MM, Kedar and Bilaskhani seem to get confused when the singers make rounded, bimanual gestures, as when the hands seem to be moving round each other. These gestures are associated with specific melodic movements, an andolan (slow oscillation) on the Ga (3) in MM and distinctive crooked (vakra) pitch movements in the other two ragas. For Sch, several ragas are wrongly predicted as Shree when she makes a direct upward hand movement, which is often distinctive of this raga [2]. Such qualitative observations suggest that the video prediction system may be classifying the hand movements in meaningful ways even when the predictions are wrong, pointing to similarities between the ways each singer gestures in different ragas.

When we look at the proportion of clips correctly identified in at least one of the two modalities in Figure 4, the most common result is to be correctly classified from the audio but incorrectly classified from the video. The small percentage of clips for which the opposite is true (c. 6.6 %) suggests there is some scope for the video information to improve the prediction accuracy achieved through audio alone, although this seems like a difficult challenge as the audio prediction accuracy is already high. Even so, we note that the 6.6 % where the audio is incorrect but video correct, the multimodal condition actually recovers about a third of this. Further, we have 2.5 % of the total set of clips (75/3021) correctly predicted only in the multimodal condition (i.e. audio and video data are combined using latent fusion). This amounts to a quarter of all clips that were wrongly classified in both audio-only and video-only conditions. It is interesting to note that of the 75 clips correctly predicted only in the multimodal condition, the most common ragas represented were Bahar (19) and MM (14), which as noted above share the same scale and many melodic movements. This study suggests that the combination of coordinated audio and gesture features can improve the classification of ragas from short clips. This approach could be extended to other musicological investigations such as the musical expression of mood or character, melodic phrase segmentation, and interpersonal coordination.

6. REFERENCES

- [1] M. Clayton, "Time, gesture and attention in a khyāl performance," *Asian Music*, vol. 38, no. 2, pp. 71–96, 2007.
- [2] L. Leante, "The lotus and the king: Imagery, gesture and meaning in a hindustani rāg," *Ethnomusicology Forum*, vol. 18, no. 2, pp. 185–206, 2009.
- [3] —, "Gesture and imagery in music performance: Perspectives from north indian classical music," in *The Routledge Companion to Music and Visual Culture*, T. Shephard and A. Leonard, Eds. Routledge, 2013, pp. 145–152.
- [4] —, "The cuckoo's song : imagery and movement in monsoon ragas," in *Monsoon feelings : a history of emotions in the rain.*, I. Rajamani, M. Pernau, and K. R. B. Schofield, Eds. New Delhi: Niyogi Books, 2018.
- [5] M. Rahaim, *Musicking Bodies: Gesture and Voice in Hindustani Music*. Wesleyan University Press, 2012.
- [6] S. Paschalidou, T. Eerola, and M. Clayton, "Voice and movement as predictors of gesture types and physical effort in virtual object interactions of classical indian singing," in *Proc. of the 3rd Int. Symposium on Movement and Computing*, 2016.
- [7] M. Clayton, K. Jakubowski, and T. Eerola, "Interpersonal entrainment in indian instrumental music performance: Synchronization and movement coordination relate to tempo, dynamics, metrical and cadential structure," *Musicae Scientiae*, vol. 23, no. 3, pp. 304–331, 2019.
- [8] L. Pearson, "Gesture and the sonic event in karnatak music," *Empirical Musicology Review*, vol. 8, no. 1, pp. 2–14, 2013.
- [9] —, "Coarticulation and gesture: An analysis of melodic movement in south indian raga performance," *Music Analysis*, vol. 35, no. 3, pp. 280–313, 2016.
- [10] G. Koduri, S. Gulati, P. Rao, and X. Serra, "Rāga recognition based on pitch distribution methods," *Journal of New Music Research*, vol. 41, no. 4, pp. 337–350, 2012.
- [11] K. K. Ganguli and P. Rao, "On the distributional representation of ragas: Experiments with allied raga pairs," vol. 1, no. 1, pp. 79–95, 2018.
- [12] Z. Duan, S. Essid, C. C. Liem, G. Richard, and G. Sharma, "Audiovisual analysis of music performances: Overview of an emerging field," *IEEE Signal Processing Magazine*, vol. 36, no. 1, pp. 63–73, 2019.
- [13] A. Bazzica, J. C. van Gemert, C. C. S. Liem, and A. Hanjalic, "Vision-based detection of acoustic timed events: a case study on clarinet note onsets," 2017. [Online]. Available: <https://arxiv.org/abs/1706.09556>
- [14] M. Clayton, J. Li, A. R. Clarke, M. Weinzierl, L. Leante, and S. Tarsitani, "Hindustani raga and singer classification using pose estimation," 2021. [Online]. Available: <https://doi.org/10.17605/OSF.IO/T5BWA>
- [15] Z. Cao, G. Hidalgo, T. Simon, S.-E. Wei, and Y. Sheikh, "Openpose: Realtime multi-person 2d pose estimation using part affinity fields," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 1, pp. 172–186, 2021.
- [16] R. Hennequin, A. Khelif, F. Voituret, and M. Moussallam, "Spleeter: a fast and efficient music source separation tool with pre-trained models," *Journal of Open Source Software*, vol. 5, p. 2154, 2020.
- [17] Y. Jadoul, B. Thompson, and B. de Boer, "Introducing parselmouth: A python interface to praat," *Journal of Phonetics*, vol. 71, pp. 1–15, 2018.
- [18] P. Boersma and D. Weenink, "Praat: doing phonetics by computer [Computer program]," <http://www.praat.org/>, 2021, version 6.1.38, retrieved May 2022.
- [19] J. Salamon, S. Gulati, and X. Serra, "A multipitch approach to tonic identification in indian classical music," in *Proc. of the 13th Int. Soc. for Music Information Retrieval Conference*, Porto, Portugal, 2012.
- [20] D. Bogdanov, N. Wack, E. Gómez, S. Gulati, P. Herrera, O. Mayor *et al.*, "Essentia: an audio analysis library for music information retrieval," in *Proc. of the 14th Int. Soc. for Music Information Retrieval Conference*, Curitiba, Brazil, 2013.
- [21] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proc. of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [22] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [23] C.-L. Liu, W.-H. Hsaio, and Y.-C. Tu, "Time series classification with multivariate convolutional neural network," *IEEE Transactions on Industrial Electronics*, vol. 66, no. 6, pp. 4788–4797, 2018.
- [24] H. I. Fawaz, B. Lucas, G. Forestier, C. Pelletier, D. F. Schmidt, J. Weber, G. I. Webb, L. Idoumghar, P. Muller, and F. Petitjean, "Inceptiontime: Finding alexnet for time series classification," *Data Mining and Knowledge Discovery*, vol. 34, no. 6, pp. 1936–1962, 2020.
- [25] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl, "Algorithms for hyper-parameter optimization," in *Advances in Neural Information Processing Systems*, 2011.

- [26] L. Biewald, “Experiment tracking with weights and biases,” 2020, software available from wandb.com. [Online]. Available: <https://www.wandb.com/>
- [27] T. Baltrušaitis, C. Ahuja, and L.-P. Morency, “Multimodal machine learning: A survey and taxonomy,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 2, pp. 423–443, 2018.
- [28] Z.-H. Zhou, *Ensemble methods: foundations and algorithms*. CRC press, 2012.
- [29] N. Hatami and R. Ebrahimpour, “Combining multiple classifiers: diversify with boosting and combining by stacking,” *Int. Journal of Computer Science and Network Security*, vol. 7, no. 1, pp. 127–131, 2007.
- [30] S. Džeroski and B. Ženko, “Is combining classifiers with stacking better than selecting the best one?” *Machine learning*, vol. 54, no. 3, pp. 255–273, 2004.
- [31] J. H. Koo, S. W. Cho, N. R. Baek, M. C. Kim, and K. R. Park, “Cnn-based multimodal human recognition in surveillance environments,” *Sensors*, vol. 18, no. 9, p. 3040, 2018.
- [32] L. Breiman, “Random forests,” *Machine learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [33] B. E. Boser, I. M. Guyon, and V. N. Vapnik, “A training algorithm for optimal margin classifiers,” in *Proc. of the 5th Annual ACM Workshop on Computational Learning Theory*, pp. 144–152.
- [34] R. E. Wright, “Logistic regression.” 1995.
- [35] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, “Scikit-learn: Machine learning in Python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [36] S. Oramas, F. Barbieri, O. Nieto, and X. Serra, “Multimodal deep learning for music genre classification,” *Transactions of the Int. Soc. for Music Information Retrieval*, vol. 1, no. 1, pp. 4–21, 2018.
- [37] S. Rao and W. van der Meer, “Miyan ki malhar,” <https://autrimncpa.wordpress.com/miyan-ki-malhar/>, accessed: May-2022.

TRACES OF GLOBALIZATION IN ONLINE MUSIC CONSUMPTION PATTERNS AND RESULTS OF RECOMMENDATION ALGORITHMS

Oleg Lesota^{1,2}

Emilia Parada-Cabaleiro^{1,2}

Stefan Brandl¹

Elisabeth Lex³

Navid Rekabsaz^{1,2}

Markus Schedl^{1,2}

¹Multimedia Mining and Search Group, Institute of Computational Perception, JKU Linz, Austria

²Human-centered AI Group, AI Lab, Linz Institute of Technology (LIT), Austria

³Graz University of Technology, Austria

markus.schedl@jku.at

ABSTRACT

Music streaming platforms allow users to enjoy music from all over the globe. Such opportunity speeds up cultural exchange between different countries, a process often associated with globalization. While such an exchange could lead to more diverse music consumption, empirical evidence on its influence on online music consumption is limited. Besides, the extent to which music recommender systems foster exchange or amplify globalization in music remains an understudied problem.

In this paper, we present findings from an empirical study to detect traces of globalization in domestic vs. foreign online music consumption. Besides, we investigate if popular recommendation algorithms, specifically ItemKNN and NeuMF, are prone to amplifying globalization processes. Our experiments on Last.fm listening data show nuanced patterns of globalization in music consumption. We observe a strong position of US music in all considered countries. In countries such as Sweden, Great Britain, or Brazil, US music shows various levels of coexistence with domestic music. We find that Finland is least influenced by US music, while greatly consuming and “exporting” domestic music. With respect to recommendation algorithms, ItemKNN tends to recommend domestic music to users of many countries, while NeuMF contributes to accelerating globalization and shifting balance towards dominance of US music on the market.

1. INTRODUCTION AND RELATED WORK

Globalization can be defined as an “*expanding cultural exchange between countries, which may imply an increasing consumption of foreign cultural goods beside the local ones*” [1]. Among others, previous research proposes two interpretations of globalization: (i) *Cultural Imperialism*

[2], i. e., the growing cultural exchange triggered by globalization is mostly profitable for certain dominant western cultures (in particular American culture) and thereby threatens to overwhelm the others; and (ii) *Glocalization* [1, 3], i. e., fostering the development of local cultures through the globalization process, by means of adaptation of global cultural forms and strengthening local identity as a counterbalancing mechanism against global influences. These interpretations imply that globalization exposes local cultures to pressure from global trends, and while some of them adapt and confront the threat, others weaken and decline.

The realizations of these interpretations in various domains have been confirmed by several studies. For instance, Crane [2] shows the dominance of US film industry in most regions of the world, while Chen and Shen [4] display the ability of some cultures to adapt and develop under the pressure of more dominant ones. Approaching globalization, most of the previous works resort to the analysis of various aggregated popular charts representing mainstream consumption trends [1, 5–7]. Specifically, Achterberg [1] show that US music has become increasingly popular in the Netherlands, Germany, and France until the late 1980s’, while starting in the 90s’, more local music has been produced. The revival of domestic music consumption is also evidenced by Bekhuis et al. [6], who show that popularity of domestic artists is positively correlated with a high sentiment of national pride. Similarly, Verboord and Brandellero [7] conduct a multilevel analysis of pop-charts’ evolution, showing that the consumption of foreign music has increased in many countries except in the US.

The mentioned studies conduct the analysis on aggregated pop-music charts, neglecting the nuances of music consumption by the consumers with less mainstream tastes. Including less-mainstream listeners in country-specific analyses is particularly important, as listeners in different countries display diverse tendencies towards listening to mainstream music [8], which also gets translated into significant differences in the *performance* of recommendation algorithms [9, 10]. However, to the best of our knowledge, little research has been dedicated to investigate the *influence* of streaming platforms on globalization processes [5]. Furthermore, the extent to which existing globalization processes might be amplified by music rec-



© O. Lesota, E. Parada-Cabaleiro, S. Brandl, E. Lex, N. Rekabsaz, and M. Schedl. Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). **Attribution:** O. Lesota, E. Parada-Cabaleiro, S. Brandl, E. Lex, N. Rekabsaz, and M. Schedl, “Traces of Globalization in Online Music Consumption Patterns and Results of Recommendation Algorithms”, in *Proc. of the 23rd Int. Society for Music Information Retrieval Conf.*, Bengaluru, India, 2022.

ommender systems has, to the best of our knowledge, not yet been investigated.

To bridge these gaps, we leverage online music listening data, allowing us to analyze globalization patterns based on actual per-user consumption of music from a wider range of countries. In addition, we study potential impacts music recommender systems may have on globalization. To this end, we strive to answer the following three research questions:

- **RQ1:** How prominent is the aspect of US “cultural imperialism” in the sphere of online music consumption? Is its influence uniform across countries?
- **RQ2:** How significant is domestic music consumption in different countries? Are there any signs of “glocalization”?
- **RQ3:** Do music recommender systems influence users’ inclination towards music coming from certain countries?

We highlight the importance of RQ3, considering recent research showing that recommender systems can amplify biases from underlying data in their output, as well as inflict additional algorithmic biases. An overview of biases recommender systems are prone to is given in [11, 12]. Among them is popularity bias, i.e., the tendency to favor popular items at the expense of niche and less popular items, often leading to dissatisfaction of users [13, 14]. We pose RQ3 to investigate possible occurrence of analogous “globalization bias” in recommender systems.

The remainder of the paper is structured as follows. Section 2 describes our methodological approach to the research questions, including data and methods used. We report and analyze our results in Section 3. Section 4 is dedicated to limitations of our study. Finally, we sum up our findings and list potential future directions in Section 5.

2. METHODS AND MATERIALS

In the following, we detail the methodology we adopt to answer the research questions and describe the dataset our analysis is based upon. We first clarify and formally define the terminology we use in the methodological description: A *listening event* (or *interaction*) is a tuple $\langle u, i \rangle$, indicating that a user u consumed an item i , which is a music track or a song in our case. Each track i has been created or performed by an artist a .¹ Each artist a originates from a country c_a . Each user u is from a country c_u . \mathcal{I}_{c_u} denotes the set of all user–item interactions made by users in country c_u . \mathcal{I}_{c_a} refers to the set of all interactions with items created by artists from country c_a . \mathcal{I}_{c_a, c_u} denotes the set of all interactions with items created by artists from country c_a , listened to by users from country c_u .

2.1 Investigating US Cultural Imperialism

To answer RQ1, we compute, for all users in a given fixed country c_u , their aggregate share of consumed music that

¹ Being aware that “artist” in the context of music can refer to different subjects, we adopt a pragmatic definition here. Owing to the type of user-generated data we investigate in our study, we consider both composers and performers as artists.

has been created by artists from each country under consideration. Since we are interested in this share for local and US music (artists), we define the following function

$$IC_{c_u \rightarrow c_a} = \frac{|\mathcal{I}_{c_a, c_u}|}{|\mathcal{I}_{c_u}|}$$

for which we consider three cases:

- *Local:* $c_a = c_u$,
- *US:* $c_a = US$, and
- *Other:* $\sum_{c_a \in C \setminus \{c_u, US\}} IC_{c_u \rightarrow c_a}$, where C is the set of all countries considered (sum over all artist countries that are not the local country nor the US).

In simple words, this formal framework can answer how much of the music all users in country c_u consume was created by artists from the same country (local consumption), was created by artists from the US, and was created by artists from other countries (neither local nor the US); all expressed in relative numbers.

2.2 Investigating Traces of Glocalization

To answer RQ2, we consider all listening events of songs by artists from the country under investigation, i.e., we fix c_a . We then define $IC_{c_u \leftarrow c_a}$, analogously to $IC_{c_u \rightarrow c_a}$ above, i.e., for a given country c_a , the share of its artists’ listening events that originate from users in each country c_u under investigation.

$$IC_{c_u \leftarrow c_a} = \frac{|\mathcal{I}_{c_a, c_u}|}{|\mathcal{I}_{c_a}|}$$

Here we are mostly interested in this share between local and foreign music consumers, and accordingly consider two cases:

- *Local:* $c_u = c_a$, and
- *Other:* $\sum_{c_u \in C \setminus \{c_a\}} IC_{c_u \leftarrow c_a}$, where C is the set of all countries considered (sum over all user countries that are not the local country).

In simple words, this formal framework can answer how much of the music created by artists from c_a is consumed by users in the same country (local consumption of local artists), and how much by users in other countries; all expressed in relative numbers.

2.3 Influence of Recommendation Algorithms on Globalization Patterns

To approach RQ3, we consider two popular recommendation algorithms, a classical ItemKNN [15] and a more recent, deep-learning-based NeuMF [16]. ItemKNN assigns a recommendation score to an item for a given user based on how similar this item is to the items already consumed by the user. Similarity is computed based on the interactions of other users. This algorithm does not learn any special representations for users and items, operating directly on the interaction matrix. On the other hand, NeuMF is a matrix factorization approach. Not only does it learn user and item embeddings, but also a dedicated scoring