
Visualizing and Understanding Self-attention based Music Tagging

Minz Won¹ Sanghyuk Chun² Xavier Serra¹

Abstract

Recently, we proposed a self-attention based music tagging model (Won et al., 2019). Different from most of the conventional deep architectures in music information retrieval, which use stacked 3×3 filters by treating music spectrograms as images, the proposed self-attention based model attempted to regard music as a temporal sequence of individual audio events. Not only the performance, but it could also facilitate better interpretability. In this paper, we mainly focus on visualizing and understanding the proposed self-attention based music tagging model.

However, the aforementioned models for MIR are yet less interpretable. All of the introduced models are using stacked 3×3 filters, which was originally designed for the image processing, on spectrogram inputs. This architecture captures spectro-temporal local features in each layer, while music is a temporal sequence of individual audio events. From this motivation, we recently proposed a new architecture design for music tagging that captures timbral local features using CNN and learns their temporal relation using self-attention mechanism (Won et al., 2019). The proposed model is not only good in its performance but also facilitates interpretable visualization. In following sections, we summarize the model architecture (Section 2), visualize learned information (Section 3), and finalize the paper by discussing future work (Section 4).

1. Introduction

As deep learning based research successfully demonstrated its versatility, interest in interpretability of deep learning models has been increased together. In the field of computer vision (CV), researchers tried to interpret the reasons and mechanisms behind the predictions of deep architectures by: mapping intermediate feature activities to the input space using deconvolution (Zeiler & Fergus, 2014), highlighting class discriminative localization map using gradients of the target class (Grad-CAM) (Selvaraju et al., 2017), and perturbing input to determine local importance (LIME) (Ribeiro et al., 2016). Motivated by the previous works, music information retrieval (MIR) researchers also endeavored to understand their deep architectures. Visualization using deconvolution was demonstrated with a genre classification model (Choi et al., 2016) and an instrument recognition model (Han et al., 2017). Especially, Choi et al. proposed an auralization method to interpret the network by listening to the reconstructed signal from deconvolved spectrograms (2016). To explain the predictions of a singing voice detection model, Mishira et al. adopted LIME (2017) and feature inversion (2018).

2. Model Architecture

The self-attention based music tagging model (Won et al., 2019) consists of two parts: front-end that captures timbral local information and back-end that models temporal relation of the captured local features. Although two different front-ends were introduced in the paper, in this work, we only use spectrogram-based CNN which is equivalent to the front-end of Pons et al. (2018). The front-end CNN consists of vertical (i.e. 86×7) and horizontal (i.e. 1×129) filters.

The back-end of the model is identical to the encoder of the Transformer (Vaswani et al., 2017). It consists of stacks of self-attention layers and was originally designed to solve natural language processing tasks. A self-attention module computes the response at a location in a sequence by attending to all locations within the same sequence. Since music composes its semantics based on the relations between components in sparse positions, we can adopt self-attention layers for music sequence modeling. Although features from the front-end are not discrete like language, self-attention has already been successfully plugged into computer vision architectures (Wang et al., 2018) which use linear features.

In this work, we used stacks of 2 multi-head attention layers with 8 attention heads which showed the best performance for $\approx 4.1s$ input sequence. Implementation details are available online¹ for reproducibility.

¹Music Technology Group, Universitat Pompeu Fabra, Barcelona, Spain ²Clova AI Research, Naver Corp., Seongnam, Korea. Correspondence to: Minz Won <minz.won@upf.edu>.

¹<https://github.com/minzwon/self-attention-music-tagging>

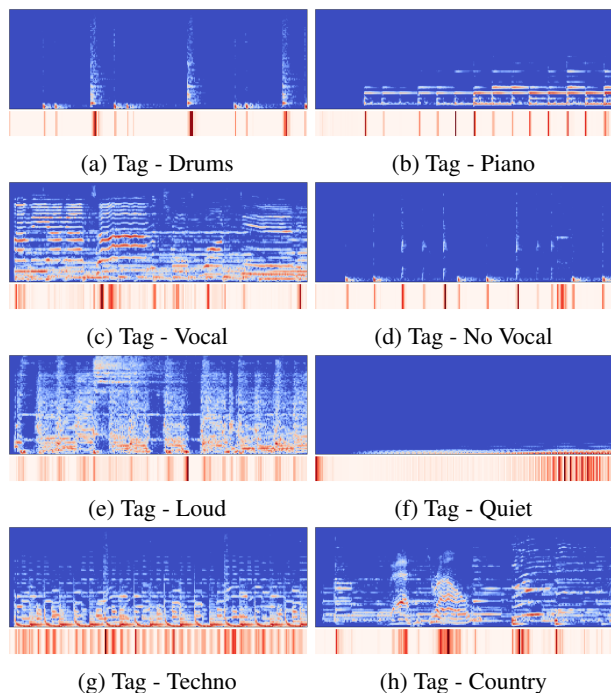


Figure 1. Attention heat maps.

3. Visualization

Attention Heat Map. To understand the behavior of the model, it is important to know which part of the audio the machine pays more attention to. To this end, we summed up attention scores from each attention head and visualized them. Figure 1 shows log mel-spectrograms and their attention heat maps. For simplification, we only visualized the attention heat map of the last attention layer. The model tends to pay more attention to more informative parts for music tagging. However, as shown in Figure 1d and 1f, attention heat maps always highlight parts with more energy although they were tagged as *no vocal* and *quiet*, respectively. Also, it is difficult to interpret the reason of tagging if the tags are related to longer-term information (Figure 1e and 1f). Attention heat maps can pinpoint where the machine pays attention, but they cannot provide reasons for the classification or tagging.

Tag-wise Contribution Heat Map. To emphasize which part of the audio is more relevant to each tag, we visualized tag-wise contribution heat maps (Figure 2). We manually changed the attention score of the last attention layer. For each time step, we manipulated the attention score as 1 and set scores of other parts as 0 so that we can see the contribution of each time bin to each tag. This tag-wise contribution heat map is inspired by the manual attention weight adjustment proposed by Lee et al. (2017). To compare the different contribution of different audio, we concatenated two spectrograms and fed them through the network. For

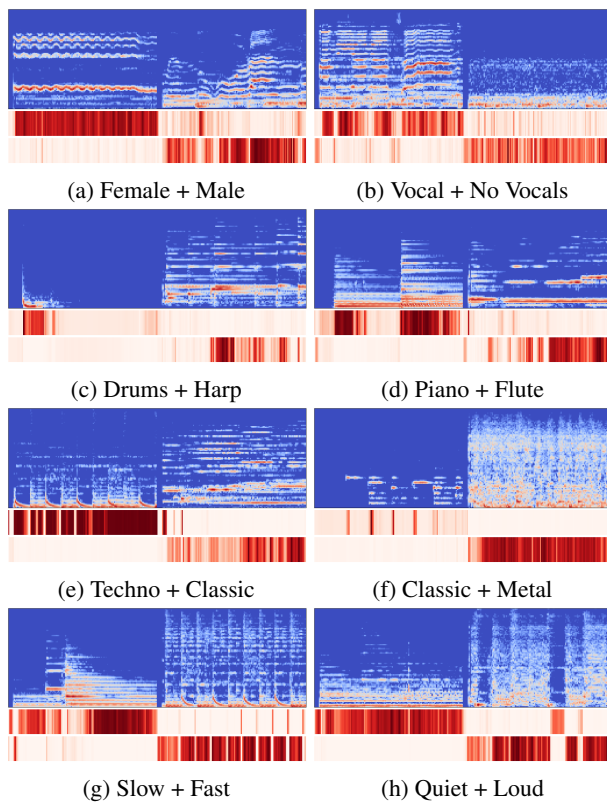


Figure 2. Tag-wise contribution heat maps.

example, Figure 2a shows a concatenated spectrogram of *female* and *male* voice. The first row heat map highlights contribution to the *female* tag and the second row indicates contribution to the *male* tag. We repeated this for instruments (2a, 2b, 2c, 2d), genres (2e, 2f), and moods/themes (2g, 2h). Different from the attention heat map, the tag-wise contribution heat map can facilitate better interpretation of the classification since it visualizes contribution to each tag. By comparing Figure 1d and 2b, we could figure out that the model pays more attention to parts with more energy but the contribution to *no vocal* tag is different from the attention heat map. Interestingly, the second half of the spectrogram in Figure 2h has a temporary silence and the contribution heat map for *quiet* shows an according short activation.

4. Future Work

We demonstrated the interpretability of our proposed model with a use case of music tagging. Since the architecture design is not task-specific, it can be applied to solve general MIR problems. However, the front-end of the proposed model is yet less interpretable. We could highlight relevant parts of audio for each tag, but still don't know which frequency band or timbral information is important for each tag. By adopting gradient-based visualization methods in the front-end, one can expect better interpretability.

Acknowledgements

This work was funded by the predoctoral grant MDM-2015-0502-17-2 from the Spanish Ministry of Economy and Competitiveness linked to the Maria de Maeztu Units of Excellence Programme (MDM-2015-0502).

References

- Choi, K., Fazekas, G., and Sandler, M. Explaining deep convolutional neural networks on music classification. *arXiv preprint arXiv:1607.02444*, 2016.
- Han, Y., Kim, J., Lee, K., Han, Y., Kim, J., and Lee, K. Deep convolutional neural networks for predominant instrument recognition in polyphonic music. *IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP)*, 25(1):208–221, 2017.
- Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.
- Lee, J., Shin, J.-H., and Kim, J.-S. Interactive visualization and manipulation of attention-based neural machine translation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 121–126, 2017.
- Mishra, S., Sturm, B. L., and Dixon, S. Local interpretable model-agnostic explanations for music content analysis. In *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, pp. 537–543, 2017.
- Mishra, S., Sturm, B. L., and Dixon, S. Understanding a deep machine listening model through feature inversion. In *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, 2018.
- Pons, J., Nieto, O., Prockup, M., Schmidt, E., Ehmann, A., and Serra, X. End-to-end learning for music audio tagging at scale. In *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, 2018.
- Ribeiro, M. T., Singh, S., and Guestrin, C. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135–1144. ACM, 2016.
- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 618–626, 2017.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017.
- Wang, X., Girshick, R., Gupta, A., and He, K. Non-local neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7794–7803, 2018.
- Won, M., Chun, S., and Serra, X. Toward interpretable music tagging with self-attention. *arXiv preprint arXiv:1906.04972*, 2019.
- Zeiler, M. D. and Fergus, R. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pp. 818–833. Springer, 2014.