

AESTHETIC MATTERS IN MUSIC PERCEPTION FOR IMAGE STYLIZATION: A EMOTION-DRIVEN MUSIC-TO-VISUAL MANIPULATION

Junjie Xu¹, Xingjiao Wu¹, Tanren Yao¹, Zihao Zhang¹, Jiayang Bei¹, Wu Wen¹, Liang He¹

¹East China Normal University, Shanghai, China

ABSTRACT

Emotional information is essential for enhancing human-computer interaction and deepening image understanding. However, while deep learning has advanced image recognition, the intuitive understanding and precise control of emotional expression in images remain challenging. Similarly, music research largely focuses on theoretical aspects, with limited exploration of its emotional dimensions and their integration with visual arts. To address these gaps, we introduce EmoMV, an emotion-driven music-to-visual manipulation method that manipulates images based on musical emotions. EmoMV combines bottom-up processing of music elements such as pitch and rhythm with top-down application of these emotions to visual aspects like color and lighting. We evaluate EmoMV using a multi-scale framework that includes image quality metrics, aesthetic assessments, and EEG measurements to capture real-time emotional responses. Our results demonstrate that EmoMV effectively translates music's emotional content into visually compelling images, advancing multimodal emotional integration and opening new avenues for creative industries and interactive technologies.

Index Terms— Multimodal Generation, Cross-modal Learning, Music Emotion Recognition, Image Aesthetics

1. INTRODUCTION

Images contain multi-level information, including low-level features (e.g., color, texture), mid-level features (e.g., shape, edges), high-level objects (e.g., people, scenes), and emotional and aesthetic content [1]. Emotional information is crucial for improving human-computer interaction, advancing creative industries, and deepening image understanding. However, despite advancements in deep learning for image recognition and classification, emotional expression in images remains underexplored. Specifically, **intuitively understanding and precisely controlling emotional expression in images is challenging**, making it a persistent research problem.

Music, a powerful medium for emotional expression, can evoke human emotional resonance. Understanding its emotional content not only deepens music comprehension but also enhances higher-level cognitive understanding. Cur-

rent research on music mainly focuses on theoretical tasks like instrument recognition and beat analysis, with limited exploration of emotional expression. While some studies use emotion-annotated datasets and classification models [2], they rely on textual annotations and fail to capture the fine-grained, multi-dimensional nature of musical emotions. Therefore, **the richness and complexity of musical emotions cannot be fully encapsulated by simple textual descriptions**. Many creative studies [3, 4, 5] combine music and visual arts to evoke emotions, and studies [6, 7] show that their synergy deepens emotional understanding. Integrating visual imagery can help extract emotional information from music, while music's emotional cues can enhance the emotional depth of images. This interdisciplinary fusion improves emotional expression in visuals and introduces new methods for emotion-driven image generation. Thus, a multimodal approach to linking music and images is valuable. While studies [8, 9, 10, 11] have explored cross-modal relationships between audio and visual information, they often focus on semantic extraction, neglecting the emotional dimension. Therefore, **correlating emotional information in music and visuals remains a significant challenge**.

Therefore, we proposed an **Emotion-driven Music-to-Visual** manipulation method, termed **EmoMV**, to manipulate an image with the condition from music. Inspired by human cognitive processes [12], we adopt a strategy that seamlessly integrates bottom-up (music-to-emotion) and top-down (emotion-to-image) stages to translate emotional information from music into visual imagery. Building upon the research of [13, 14], which demonstrates that emotional responses are elicited by the fundamental structural elements of music such as pitch, rhythm, and chord changes we begin with these basic theoretical components. By incorporating supplementary information from visual arts, we progressively map these foundational elements to higher-level emotional expression (up-bottom ↓). While, learning from the work [15] which validates aesthetic attributes can be related to emotional expression, we further embody emotions into the aesthetic-related low-level dimensions like light, exposure, and color (bottom-up ↑).

A key challenge in evaluating emotion-driven manipulations is the subjectivity of both aesthetic and emotional experiences. To address this, we propose a multi-scale evalua-

tion approach combining image quality assessments, aesthetic evaluations, and EEG measurements. EEG tracks real-time emotional responses to images, providing a more objective and immersive understanding of how well the images capture the music’s emotional content. This dual evaluation offers a comprehensive measure of our methods’ effectiveness in translating music’s emotional qualities into visual output.

The Contributions of this paper are as follows:

- We develop a two-stage framework called EmoMV, which enables emotion-driven image manipulation. EmoMV leverages bottom-up and top-down cognitive strategies to construct emotional representations of music and images, in line with cognitive science principles.
- We introduce a Mus-Vis Textual Alignment module for aligning musical description and visual description with emotional expression. While Emotion-aware Aesthetic Image Refinement is proposed for mapping emotional information into visual elements like lighting, and exposure, and achieving the harmony of the whole image.
- We propose a multi-scale evaluation framework and demonstrate the effectiveness of EmoMV through extensive experiments on a 38k music-image pair dataset that we collected from online sources. Our results show that EmoMV can positively impact emotional well-being, particularly in applications such as art therapy.

2. RELATED WORK

Image Emotion. Recent studies have explored emotional content in images by designing and extracting visual features at different levels. Yanulevskaya [16] proposed an emotional classification method for artworks based on low-level features. The influence of composition, color contrast, and texture on emotional expression is discussed in [17], where features are designed according to artistic principles. Traditional methods, however, failed to capture all factors influencing human emotions. More recent approaches focus on extracting semantic features but often overlook higher-level abstract elements. Yang [18] used object detection and attention mechanisms for emotion recognition, while [19] introduced a network to learn emotional correlations from visual stimuli. Despite these efforts, the abstract nature of emotion still poses challenges, requiring the integration of auxiliary information for more accurate visual emotion recognition.

Music Emotion. Music emotion research [20], an interdisciplinary field combining music theory, cognitive science, and AI, has gained significant attention. Researchers have explored emotion recognition [21], classification [22], and generation [23, 24]. A key challenge is effectively representing

emotions. Traditional methods use discrete emotion models, identifying emotions like happiness or sadness. However, Chaturvedi [25] highlights the multidimensional nature of emotional experiences, advocating for more nuanced frameworks beyond simple annotations. Relying on subjective feedback or textual descriptions alone fails to capture the full complexity of musical emotions.

Music-image Alignment. The alignment of music and visual [26, 27, 28] to evoke and enhance emotional responses has been a key focus in the field of deep learning. Music has the ability to evoke complex emotional states, and when combined with visual, it can provide a richer emotional context [6, 7, 29]. This combination not only stimulates the emotional centers of the brain but also enhances the expressiveness of visual content. While recent studies explore music-image fusion, most [8, 9, 10, 11] focus on semantic content, often neglecting the emotional dimension. Effectively correlating emotional information across these modalities remains challenging, as emotional alignment requires capturing subtle nuances, not just explicit features [30, 31]. Additionally, emotions vary across music genres [32] and visual styles.

EEG-based Emotional Perception. Emotional perception in the audio-visual modality has become a key research focus. Visual and auditory stimuli each evoke distinct emotions, with facial expressions and emotional images triggering feelings like happiness or fear [33], while music and vocal tone influence emotional perception [34]. When these cues align, emotional perception is consistent; however, mismatched cues can lead to fluctuating or ambiguous responses [35]. EEG has become a valuable tool for tracking emotional fluctuations induced by audiovisual stimuli, providing an objective measure of emotional shifts, particularly in complex multimodal contexts [36]. Therefore, we propose an EEG-based Multi-scale evaluation framework for evaluating our method.

3. METHODOLOGY

In this section, we detail EmoMV, a two-stage approach comprising a *Mus-Vis Textual Alignment* and a *Emotion-aware Aesthetic Image Refinement*. The former correlates the musical elements with emotional expression, while the latter externalizes emotional information into the aesthetic attributes of the image.

3.1. Task Definition

Given a music piece M and an image I , the goal is to design a method F that generates an emotionally relevant image I_{final} while preserving the semantic content of I : $I_{final} = F(I, M)$.

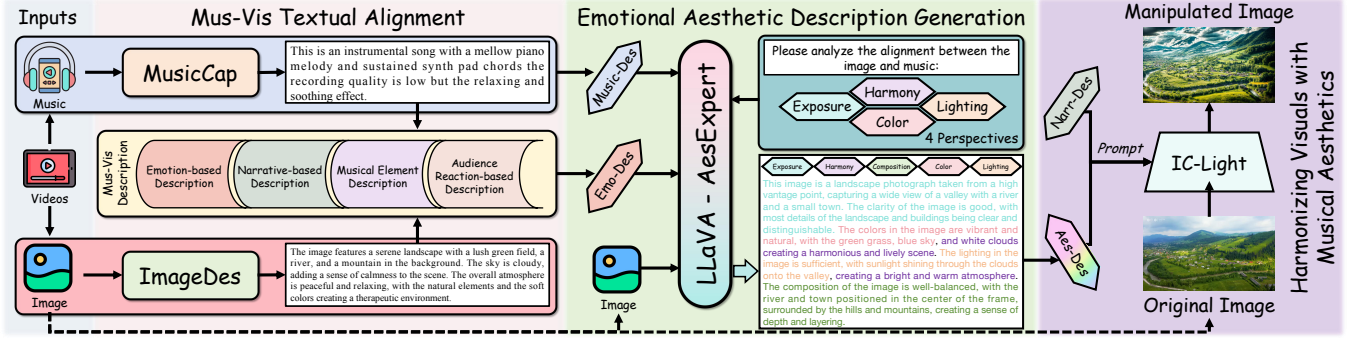


Fig. 1: Overview of our framework Emotion-driven Music-to-Visual manipulation method, termed as **EmoMV**. **EmoMV** is composed of *Mus-Vis Textual Alignment*, *Emotional Aesthetic Description via Lighting*, and *Harmonizing Visuals with Musical Aesthetics*, facilitating the extraction of emotional information and exhibit on the generated image.

3.2. Mus-Vis Textual Alignment (MVTA)

To extract emotional expression from music, we elevate our understanding of the fundamental elements of music to the level of emotional expression. In this process, we propose a novel emotion alignment method that incorporates image-based assistance, leveraging the specific information of the image being edited to enhance the emotional expression of music in relation to the image. We begin by extracting the basic elements of music, drawing inspiration from [37]. We introduce an encoder-decoder framework to generate musical element descriptions, using the powerful audio encoder HT-SAT [38] for music encoding and BART for decoding. The model is trained on the MusicCaps [39] and MusicBench [40] datasets. Once the musical element descriptions are obtained, we integrate them with the image input, employing the robust multimodal model LLaVA to generate four-dimensional descriptions (emotional, narrative, musical, and audience reaction). This allows us to map the music element descriptions and the image input into emotion-related descriptions:

$$Description = LLaVA\{[HT - SAT + BART](M), I\} \quad (1)$$

3.3. Emotion-aware Aesthetic Image Refinement (EAIR)

The overall emotional expression of an image originates from pixel-level changes. We map the abstract emotional expression in music to fundamental aesthetic attributes, applying a top-down approach to translate musical emotions into the core elements of the image. This process consists of two main parts:

Emotional Aesthetic Description Generation: By utilizing music descriptions, emotional annotations, and the image itself, we fine-tune the AesExpert model to embed emotional descriptions into the aesthetic attributes of the image. This step allows us to capture how emotional expressions manifest through various visual attributes of the image.

$$Aes - Des = AesExpert(I, Emo - Des, Music - des) \quad (2)$$

Harmonizing Visuals with Musical Aesthetics: We use the IC-Light model, along with our image emotion attribute descriptions, to edit the image. Additionally, we incorporate Narrative descriptions for semantic completion, ensuring semantic consistency before and after the image editing process.

$$I_{final} = IC - Light(I, Aes - Des, Narr - Des) \quad (3)$$

Table 1: Comparison of Image Quality Metrics for *HealSoul-1k* and *EmoMV-1k* Datasets. Better results are highlighted in pink.

Metric	HealSoul-1k	EmoMV-1k
Sharpness	(Mean Std) 354.48 ± 440.19	656.58 ± 376.85
	(Min, Max) (1.63, 4292.79)	(3.42, 1900.68)
Contrast	(Mean Std) 0.97 ± 0.07	0.99 ± 0.04
	(Min, Max) (0.40, 1.00)	(0.36, 1.00)
Colorfulness	(Mean Std) 143.74 ± 29.33	142.73 ± 22.47
	(Min, Max) (35.48, 186.77)	(60.94, 184.75)
BRISQUE	(Mean Std) 37.13 ± 16.75	28.07 ± 13.31
	(Min, Max) (1.94, 106.20)	(-0.97, 71.91)

4. EXPERIMENTS

4.1. Datasets & Experiment Settings

We collected a dataset of 3.8k music-image pairs by searching for healing-related keywords on online resource and conducted experiments on this dataset to demonstrate the effectiveness of our method. We propose a multi-dimensional evaluation framework that includes metrics from the perspectives of image aesthetics, image quality, and human feedback. Specifically, we assess the following aspects:



Fig. 2: A qualitative analysis is conducted by contrasting the original image from HealSoul and the generated image through EmoMV. With the paired healing music, EmoMV generates images with more light, better composition.

Table 2: Aesthetic comparison of images derived from VEGAS [41], VGGSound [42], HealSoul-1k, and EmoMV-1k. The best results are highlighted in pink, while suboptimal results are marked in green.

Datasets		VEGAS[41]	VGGSound[42]	HealSoul-1k	EmoMV-1k
VILA	mean \uparrow	0.1891	0.3232	0.5694	0.6304
	variance \downarrow	0.0042	0.0152	0.0138	0.0084

VILA Score: In image synthesis, evaluating aesthetic quality is critical. To reduce subjective bias, we adopt the VILA framework [43], an advanced method for objective and reliable aesthetic assessment.

Image Quality Assessment: We assess image quality using four key metrics: **Sharpness**, which quantifies edge clarity and detail; **Contrast**, evaluating luminance and color differences to enhance visual depth; **Colorfulness**, measuring the richness and vibrancy of hues; and **BRISQUE** [44], a no-reference metric that evaluates perceptual quality using natural scene statistics. These metrics collectively provide a comprehensive evaluation of visual quality.

Prefrontal EEG Analysis: We analyze prefrontal EEG signals to assess the emotional and cognitive impact of generated images. The detected brainwave activity is closely linked to users' emotional states and cognitive engagement, offering real-time, objective insights that complement traditional quality assessments.

4.2. Qualitative Analysis

We conducted a qualitative analysis to assess the impact of *EmoMV* on the emotional and aesthetic qualities of generated images. We selected 10 images from the *HealSoul* dataset and compared them with *EmoMV*-generated images using corresponding music inputs. As shown in Fig. 2, the *HealSoul* dataset, consisting of 38k natural landscape images, conveys a bright, warm aesthetic that evokes a soothing, therapeutic atmosphere. In contrast, *EmoMV*-generated images, enhanced with emotional cues from music, exhibit brighter lighting, richer stylistic expression, and improved visual quality, resulting in more vivid, emotionally resonant visuals.

Table 3: EEG signal comparison between the *HealSoul-1k* and *EmoMV-1k* datasets. Superior results are highlighted in pink.

Frequency bands		<i>HealSoul-1k</i>	<i>EmoMV-1k</i>
Alpha \downarrow	Mean	27306.10	27905.21
	Events	5	4
Beta \uparrow	Mean	19211.42	27996.47
	Events	8	21
Gamma \uparrow	Mean	7691.83	9028.21
	Events	5	8

4.3. Quantitative Analysis

We designed a multi-scale evaluation system that includes image quality assessment (objective), aesthetic large model evaluation (semi-objective and semi-subjective), and EEG monitoring (subjective). As shown in Tab. 1 & 2, the *HealSoul* dataset demonstrates higher aesthetic quality, while images generated by *EmoMV* exhibit superior image quality and aesthetic performance compared to *HealSoul*. Through emotional guidance via music, *EmoMV* enhances both the quality and aesthetic appeal of the images.

As demonstrated in [45], alpha waves are negatively correlated with positive emotions, whereas beta and gamma waves are positively correlated with good emotions. The “mean” amplitude of these waves reflects the overall emotional state, while “events” signify sharp fluctuations in brainwave activity, indicating emotional tension. To evaluate the effectiveness of *EmoMV*, we conducted a study with 15 participants who were exposed to 100 synchronized



Fig. 3: Ablation study for the effect of prompt input into the iclight, we put images from *HealSoul* and generated from EmoMV branches without EAIR (w/o EAIR) or without MVTA (w/o MVTA).

image-music pairs while their EEG brainwave responses were recorded (see Table 3). The results indicate that EmoMV significantly enhances positive emotional responses and reduces negative emotional responses compared to baseline methods. Specifically, EmoMV-generated images elicited more frequent increases in positive emotions and fewer instances of negative emotions, demonstrating superior emotional expression and audiovisual consistency. These findings validate EmoMV’s ability to generate images with coherent and consistent emotional dimensions.

4.4. Ablation Study

To evaluate the roles of the *MVTA* and *EAIR* modules in EmoMV, we conducted ablation studies. Removing the *MVTA* module blocks the extraction of emotion-related and musical information, causing the EADL module to process only image inputs and pass Aes-Des directly to IC-Light. Similarly, excluding the *EAIR* module by eliminating Aes-Des from EADG disrupts the Emotional Aesthetic Description Generations ability to transform information. As shown in Table 4, EmoMV-1k achieves superior image quality, particularly in Sharpness and BRISQUE metrics. Figure 3 illustrates that both the “w/o EAIR” and “w/o MVTA” variants exhibit reduced consistency with the original image to varying degrees. Specifically, “w/o MVTA” lacks narrative descriptions, creating a semantic gap, while “w/o EAIR” produces chaotic images that fail to convey harmony or emotional coherence. These results demonstrate the essential contributions of both modules to maintaining image quality and emotional fidelity in EmoMV.

4.5. Caes Study

This case study examines how different musical emotions influence the aesthetic stylization of images generated by EmoMV. We present two scenarios: (1) the same image with

Table 4: Comparison of the mean value of Image Quality Metrics for *w/o EAIR*, *w/o MVTA*, and *EmoMV-1k* Datasets. The best results are highlighted in pink .

Method	Sharpness	Contrast	Colorfulness	BRISQUE
<i>w/o EAIR</i>	441	1.00	147.19	36.03
<i>w/o MVTA</i>	558.97	0.99	138.39	31.49
<i>EmoMV-1k</i>	656.58	0.99	142.73	28.07

different music, and (2) the same music with different images, as illustrated in Fig. 4. In the first scenario, four music tracks with varying emotional tones are applied to a single image. Despite using the same visual content, the resulting images exhibit distinct stylistic variations based on the emotional cues in the music. For example, gentle music creates warm tones, soft lighting, and smooth transitions, evoking a serene mood, while soothing music produces cooler tones and emphasizes natural elements, maintaining a tranquil atmosphere. Other tracks, with livelier or more melancholic moods, further highlight how music’s emotional properties can shape lighting, composition, and color tones, transforming the aesthetic style of the image. The second scenario applies the same soothing music to four different scene images. Here, the music remains constant, but the generated visuals vary significantly depending on the scene. Despite the same soothing, tranquil music, the visual outcomes adapt to the context of each scene. In natural landscapes, the images maintain calm, cool tones that align with the peaceful nature of the music, while urban or abstract scenes tend to reflect subtler lighting adjustments and compositional changes that preserve the serene emotional tone of the music. These variations illustrate how EmoMV can modify visual outputs based on both the emotional input from music and the inherent characteristics of the visual scene, enhancing the emotional depth



Fig. 4: Case Study of EmoMV for Emotion-Driven Music-to-Image Generation. Left: Visual outputs generated from the same image with four different musical inputs, illustrating how variations in musical emotion (e.g., tempo, rhythm, mood) influence the aesthetic styling and emotional tone of the generated images using the EmoMV method. Right: Visual outputs generated from the same musical input but with four different scene images, demonstrating how varying visual contexts affect the interpretation of the music's emotional cues and alter the visual aesthetic representation in the EmoMV framework.

and aesthetic appeal of the images. This case study demonstrates EmoMV's flexibility in using music-driven emotional cues to create context-sensitive and emotion-driven visual stylizations, whether through varying music or diverse images.

5. CONCLUSION

In this paper, we present *EmoMV*, a novel two-stage framework for emotion-driven image manipulation that integrates bottom-up and top-down cognitive strategies to construct emotional representations from music and visuals. We introduce the Mus-Vis Textual Alignment module to harmonize musical and visual descriptions with their emotional expressions and the Emotion-aware Aesthetic Image Refinement module to map these emotions into visual elements such as lighting and exposure, ensuring cohesive image harmony. Additionally, we develop a comprehensive multi-scale evaluation framework and validate EmoMV's effectiveness through extensive experiments on a 38,000 music-image pair dataset, demonstrating its ability to translate musical emotions into compelling visuals and positively impact emotional well-being, particularly in applications like art therapy. These contributions advance the field of multimodal emotional integration, offering new avenues for creative industries and enhancing human-computer interaction.

References

- [1] D. Joshi, R. Datta, E. Fedorovskaya, Q.-T. Luong, J. Z. Wang, J. Li, and J. Luo, "Aesthetics and emotions in images," *IEEE Signal Processing Magazine*, vol. 28, no. 5, pp. 94–115, 2011.
- [2] G. Yoo, S. Hong, and H. Kim, "Emotion recognition and multi-class classification in music with mfcc and machine learning," *International Journal on Advanced Science, Engineering & Information Technology*, vol. 14, no. 3, 2024.
- [3] A. Clemente, M. T. Pearce, and M. Nadal, "Musical aesthetic sensitivity," *Psychology of Aesthetics, Creativity, and the Arts*, vol. 16, no. 1, p. 58, 2022.
- [4] N. Djalalova, "Piano performance as a factor that activates students' musical and aesthetic world views and develops musical culture," *Science and innovation*, vol. 2, no. B4, pp. 339–342, 2023.
- [5] A. M. Belfi, A. Kasdan, J. Rowland, E. A. Vessel, G. G. Starr, and D. Poeppel, "Rapid timing of musical aesthetic judgments," *Journal of Experimental Psychology: General*, vol. 147, no. 10, p. 1531, 2018.
- [6] F. Talamini, G. Eller, J. Vigl, and M. Zentner, "Musical emotions affect memory for emotional pictures," *Scientific reports*, vol. 12, no. 1, p. 10636, 2022.
- [7] M. G. Boltz, B. Ebendorf, and B. Field, "Audiovisual interactions: The impact of visual information on music perception and memory," *Music Perception*, vol. 27, no. 1, pp. 43–59, 2009.
- [8] K. Sung-Bin, A. Senocak, H. Ha, A. Owens, and T.-H. Oh, "Sound to visual scene generation by audio-to-visual latent alignment," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 6430–6440.
- [9] C.-C. Lee, W.-Y. Lin, Y.-T. Shih, P.-Y. Kuo, and L. Su, "Crossing you in style: Cross-modal style transfer from music to visual arts," in *ACM International Conference on Multimedia (ACMMM)*, 2020, pp. 3219–3227.
- [10] C.-H. Wan, S.-P. Chuang, and H.-Y. Lee, "Towards audio to scene image synthesis using generative adversarial network," in *ICASSP 2019-2019 IEEE International*

Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2019, pp. 496–500.

- [11] S. H. Lee, W. Roh, W. Byeon, S. H. Yoon, C. Kim, J. Kim, and S. Kim, “Sound-guided semantic image manipulation,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 3377–3386.
- [12] R. L. Gregory, “The intelligent eye.” 1970.
- [13] D. Huron, *Sweet anticipation: Music and the psychology of expectation*. MIT press, 2008.
- [14] C. L. Krumhansl, “An exploratory study of musical emotions and psychophysiology,” *Canadian Journal of Experimental Psychology/Revue canadienne de psychologie expérimentale*, vol. 51, no. 4, p. 336, 1997.
- [15] L. Xiao, X. Wu, J. Xu, W. Li, C. Jin, and L. He, “Atlantis: Aesthetic-oriented multiple granularities fusion network for joint multimodal aspect-based sentiment analysis,” *Information Fusion*, vol. 106, p. 102304, 2024.
- [16] V. Yanulevskaya, J. C. van Gemert, K. Roth, A.-K. Herbold, N. Sebe, and J.-M. Geusebroek, “Emotional valence categorization using holistic image features,” in *2008 15th IEEE international conference on Image Processing*. IEEE, 2008, pp. 101–104.
- [17] J. Machajdik and A. Hanbury, “Affective image classification using features inspired by psychology and art theory,” in *Proceedings of the 18th ACM international conference on Multimedia*, 2010, pp. 83–92.
- [18] J. Yang, D. She, M. Sun, M.-M. Cheng, P. L. Rosin, and L. Wang, “Visual sentiment prediction based on automatic discovery of affective regions,” *IEEE Transactions on Multimedia*, vol. 20, no. 9, pp. 2513–2525, 2018.
- [19] J. Yang, J. Li, X. Wang, Y. Ding, and X. Gao, “Stimuli-aware visual emotion analysis,” *IEEE Transactions on Image Processing*, vol. 30, pp. 7432–7445, 2021.
- [20] Y. Agrawal, R. G. R. Shanker, and V. Alluri, “Transformer-based approach towards music emotion recognition from lyrics,” in *European conference on information retrieval*. Springer, 2021, pp. 167–175.
- [21] R. Sarkar, S. Choudhury, S. Dutta, A. Roy, and S. K. Saha, “Recognition of emotion in music based on deep convolutional neural network,” *Multimedia Tools and Applications*, vol. 79, no. 1, pp. 765–783, 2020.
- [22] R. Panda, R. Malheiro, and R. P. Paiva, “Audio features for music emotion recognition: a survey,” *IEEE Transactions on Affective Computing*, vol. 14, no. 1, pp. 68–88, 2020.
- [23] M. W. Lam, Q. Tian, T. Li, Z. Yin, S. Feng, M. Tu, Y. Ji, R. Xia, M. Ma, X. Song, *et al.*, “Efficient neural music generation,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [24] F. Schneider, O. Kamal, Z. Jin, and B. Schölkopf, “Moûsai: Efficient text-to-music diffusion models,” in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2024, pp. 8050–8068.
- [25] V. Chaturvedi, A. B. Kaur, V. Varshney, A. Garg, G. S. Chhabra, and M. Kumar, “Music mood and human emotion recognition based on physiological signals: a systematic review,” *Multimedia Systems*, vol. 28, no. 1, pp. 21–44, 2022.
- [26] Y. X. Chen, Yanbei and et al, “Distilling audio-visual knowledge by compositional contrastive learning,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 7016–7025.
- [27] H.-H. Wu, P. Seetharaman, K. Kumar, and J. P. Bello, “Wav2clip: Learning robust audio representations from clip,” in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2022, pp. 4563–4567.
- [28] A. Guzhov, F. Raue, J. Hees, and A. Dengel, “Audioclip: Extending clip to image, text and audio,” in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2022, pp. 976–980.
- [29] B. Ebendorf, “The impact of visual stimuli on music perception,” Ph.D. dissertation, 2007.
- [30] T. B. Janzen, M. I. A. Shirawi, S. Rotzinger, and et al., “A pilot study investigating the effect of music-based intervention on depression and anhedonia,” *Frontiers in Psychology*, 2019.
- [31] X. Lv, Y. Wang, Y. Zhang, S. Ma, J. Liu, K. Ye, Y. Wu, V. Voon, and B. Sun, “Auditory entrainment coordinates cortical-bnst-nac triple time locking to alleviate the depressive disorder,” *Cell Reports*, 2024.
- [32] E. Jostrup, M. Nyström, E. Claesdotter-Knutsson, P. Tallberg, P. Gustafsson, O. Paulander, and G. Söderlund, “Effects of stochastic vestibular stimulation on cognitive performance in children with adhd,” *Experimental Brain Research*, vol. 241, no. 11, pp. 2693–2703, 2023.
- [33] P. J. Lang, M. M. Bradley, B. N. Cuthbert, *et al.*, “International affective picture system (iaps): Technical manual and affective ratings,” *NIMH Center for the Study of Emotion and Attention*, vol. 1, no. 39-58, p. 3, 1997.

- [34] P. N. Juslin and P. Laukka, "Communication of emotions in vocal expression and music performance: Different channels, same code?" *Psychological bulletin*, vol. 129, no. 5, p. 770, 2003.
- [35] V. I. Müller, E. C. Cieslik, B. I. Turetsky, and S. B. Eickhoff, "Crossmodal interactions in audiovisual emotion processing," *Neuroimage*, vol. 60, no. 1, pp. 553–561, 2012.
- [36] K. Kamble and J. Sengupta, "A comprehensive survey on emotion recognition based on electroencephalograph (eeg) signals," *Multimedia Tools and Applications*, vol. 82, no. 18, pp. 27 269–27 304, 2023.
- [37] X. Mei, C. Meng, H. Liu, Q. Kong, T. Ko, C. Zhao, M. D. Plumbley, Y. Zou, and W. Wang, "Wavcaps: A chatgpt-assisted weakly-labelled audio captioning dataset for audio-language multimodal research," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.
- [38] K. Chen, X. Du, B. Zhu, Z. Ma, T. Berg-Kirkpatrick, and S. Dubnov, "Hts-at: A hierarchical token-semantic audio transformer for sound classification and detection," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 646–650.
- [39] A. Agostinelli, T. I. Denk, Z. Borsos, J. Engel, M. Verzetti, A. Caillon, Q. Huang, A. Jansen, A. Roberts, M. Tagliasacchi, *et al.*, "Musiclm: Generating music from text," *arXiv preprint arXiv:2301.11325*, 2023.
- [40] J. Melechovsky, Z. Guo, D. Ghosal, N. Majumder, D. Herremans, and S. Poria, "Mustango: Toward controllable text-to-music generation," in *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, 2024, pp. 8293–8316.
- [41] Y. Zhou, Z. Wang, C. Fang, T. Bui, and T. L. Berg, "Visual to sound: Generating natural sound for videos in the wild," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3550–3558.
- [42] H. Chen, W. Xie, A. Vedaldi, and A. Zisserman, "Vg-sound: A large-scale audio-visual dataset," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2020, pp. 721–725.
- [43] J. Ke, K. Ye, J. Yu, Y. Wu, P. Milanfar, and F. Yang, "Vila: Learning image aesthetics from user comments with vision-language pretraining," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 10 041–10 051.
- [44] A. Mittal, R. Soundararajan, and A. C. Bovik, "Making a completely blind image quality analyzer," *IEEE Signal processing letters*, vol. 20, no. 3, pp. 209–212, 2012.
- [45] W.-L. Zheng and B.-L. Lu, "Investigating critical frequency bands and channels for eeg-based emotion recognition with deep neural networks," *IEEE Transactions on autonomous mental development*, vol. 7, no. 3, pp. 162–175, 2015.