



ISSN: 2230-9926

Available online at <http://www.journalijdr.com>

# IJDR

International Journal of Development Research

Vol. 16, Issue, 05, pp. 70422-70426, May, 2026

<https://doi.org/10.37118/ijdr.30821.05.2026>



RESEARCH ARTICLE

OPEN ACCESS

## WISSH – WALLPAPER INNOVATOR BY SYNTHESIZING SONG HALLMARK

\*Prof. Lalita Panika, Amit Kumar Jha, Raj Sahu, Aryan Prithwani and Siddhant Shukla

Bhilai Institute of Technology, Department of Computer Science Engineering Kendri, Raipur, India

### ARTICLE INFO

#### Article History:

Received 11<sup>th</sup> February, 2026  
Received in revised form  
26<sup>th</sup> March, 2026  
Accepted 17<sup>th</sup> April, 2026  
Published online 25<sup>th</sup> May, 2026

#### Key Words:

Social media, Student Learning, Learning  
Motivation, Education management, Human  
Resource management.

\*Corresponding author: Prof. Lalita Panika

### ABSTRACT

This paper presents an intelligent system that transforms audio music into visually expressive animated wallpapers. Leveraging deep learning-based mood recognition and advanced image generation via Stable Diffusion, the project combines acoustic feature extraction with a culturally aware mood classifier to interpret musical emotions. These emotions are converted into stylistic prompts that guide the creation of aesthetic, animated wallpapers. The system integrates PyTorch for mood detection, HuggingFace Transformers for audio embeddings, and Diffusers for visual generation. The platform offers both single and batch processing, enabling immersive visual experiences from sound input. By bridging the gap between auditory emotion and visual storytelling, this project opens a pathway for more immersive multimedia experiences. The paper explores both the algorithmic design and user-centric features of the system, demonstrating potential applications in art therapy, music streaming platforms, and personalized device aesthetics.

Copyright ©2026, Prof. Lalita Panika et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Citation: Prof. Lalita Panika, Amit Kumar Jha, Raj Sahu, Aryan Prithwani and Siddhant Shukla, 2026. "WISSH – Wallpaper Innovator by Synthesizing Song Hallmark". *International Journal of Development Research*, 16, (05), 70422-70426.

## INTRODUCTION

The emotional power of music has long been recognized across cultures, serving as a universal language that conveys feelings, memories, and atmospheres. In recent years, researchers and developers have explored the intersection of sound and visuals, aiming to create more immersive and emotionally engaging multimedia experiences. Translating audio based moods into visual art not only enriches aesthetic expression but also opens novel opportunities in human-computer interaction, media personalization, and therapeutic applications. This paper introduces a next-generation AI-powered application that analyzes the emotional content of music and generates corresponding animated wallpapers. The system bridges the gap between audio perception and generative visuals by leveraging modern machine learning and diffusion models. Designed for ease of access and broad usability, the solution identifies applicable funding agency here. If none, delete this. is deployed using Streamlit as the web interface framework. The core neural network models are implemented in PyTorch, which provides robust support for deep learning, while Stable Diffusion — a state-of-the-art generative image model — is used for producing high-resolution, stylized visuals that align with the music's emotional tone. Traditional approaches to interpreting musical moods have been highly subjective, often relying on expert-curated playlists or listener ratings. While useful, such methods lack scalability and precision. With the advancement of deep learning, especially in affective computing and audio signal processing, it has become feasible to model emotional

responses computationally. In this work, both categorical mood classes (such as happy, sad, angry, relaxed) and continuous affective dimensions (such as arousal and valence) are used. This dual approach enables a nuanced understanding of the music's emotional landscape, which in turn enhances the quality and relevance of the generated visuals. The key innovation in this system lies in its dynamic synthesis of visual content that responds to a song's intrinsic features — including rhythm, tempo, instrumentation, and expressive cues. Instead of generating static wallpapers, the system produces animated scenes that evolve over time, mimicking the progression of the song itself. For example, a calm, ambient track may result in soft, pastel-toned flowing visuals, while a fast-paced, energetic song could generate bold colors, sharp transitions, and high-motion imagery. Further distinguishing this project is the incorporation of cultural context and musical genre into the generation process. By training the model on a culturally diverse dataset of music and mood mappings, the system learns to respect genre-specific emotional nuances, such as the melancholic tone of classical Indian ragas or the uplifting energy of Latin pop. This cultural sensitivity adds an extra layer of personalization and artistic depth. From a technical perspective, the application pipeline begins with the extraction of audio features using tools such as Librosa, including MFCCs, chroma features, tempo, and spectral contrast. These features are fed into a pretrained emotion recognition model based on Transformer or LSTM architectures. The predicted mood is then mapped to a set of textual prompts that guide the Stable Diffusion image synthesis engine. The prompts are crafted to reflect not only

moodkeywords but also style, color palette, and animation cues. In addition to enhancing aesthetic experiences, this project holds potential in several domains: mood-aware ambient displays for smart homes, emotionally intelligent user interfaces, and even music therapy tools that provide visual feedback for emotional healing. By transforming auditory experiences into rich, generative visuals, the system redefines how humans can interact with both music and technology. In conclusion, the proposed application exemplifies a meaningful integration of artificial intelligence, multimedia art, and affective science. It enables users to experience music in a multisensory manner, deepening their emotional connection and offering an innovative platform for digital expression. The future scope includes real-time streaming support, mobile deployment, and personalization based on user mood profiles.

**Related Work:** Prior studies in Music Emotion Recognition (MER) have laid the groundwork for understanding the emotional content of audio through computational methods. Early techniques predominantly relied on hand-crafted features such as Mel-Frequency Cepstral Coefficients (MFCCs), tempo, zero-crossing rate, chroma features, and spectral contrast. These features capture various timbral, rhythmic, and harmonic properties of music. Such features were commonly paired with traditional machine learning classifiers, including Support Vector Machines (SVMs), Random Forests, and K-Nearest Neighbors (KNN). While these methods demonstrated reasonable performance in emotion classification tasks, they often lacked the capacity to model complex temporal dependencies inherent in music signals. With the rise of deep learning, the field has seen a significant transformation. Architectures such as Convolutional Neural Networks (CNNs) and Convolutional Recurrent Neural Networks (CRNNs) have been employed to automatically learn hierarchical representations of audio data. For instance, Kim et al. (2020) demonstrated the efficacy of CRNNs in capturing both local spectral patterns and long-term temporal structures, thereby improving accuracy in music emotion recognition tasks. Recurrent models like LSTM (Long Short-Term Memory) networks and Bidirectional GRUs have also shown success in modeling the sequential nature of audio features, allowing for more nuanced mood predictions over time. In parallel, there have been remarkable advancements in the domain of generative image modeling, especially with the emergence of diffusion models. These models, particularly Stable Diffusion, have set new standards for text-to-image synthesis by producing high-resolution, semantically rich, and stylistically diverse visuals. Stable Diffusion utilizes a latent diffusion approach, where the model learns to generate image representations from denoised latent spaces, guided by textual prompts. This allows for unprecedented control over the content and style of the generated image, making it a powerful tool for creative applications.

Despite these advancements, integrating audio-based emotion recognition with generative image synthesis—especially in a dynamic, real-time, and personalized manner—remains an unexplored area. Most previous works in this space have focused on generating static visualizations from music, such as abstract art or reactive visual equalizers. Some experimental efforts have utilized Generative Adversarial Networks (GANs) for music-to-image translation. However, these approaches often suffer from limitations in image quality, interpretability, and controllability. Moreover, they typically do not incorporate emotion modeling as a core component, nor do they offer adaptive or animated outputs based on temporal emotional changes in the music. A few interactive music visualizers attempt to reflect audio features such as beat or pitch, but they rarely incorporate emotional dimensions like valence and arousal, nor do they synthesize semantically meaningful or aesthetically stylized scenes. Additionally, most of these systems operate in a generic or rule-based manner, lacking personalization, cultural context, or narrative depth. In contrast, diffusion models—especially when combined with emotion-aware conditioning—present a novel pathway for translating the affective properties of music into expressive, personalized visuals. Our approach leverages this capability by using a deep learning-based emotion inference pipeline to interpret the mood of an audio track and then

generating corresponding visual prompts. These prompts are not static but evolve over time, guiding the image generation process to produce animated, mood-adaptive wallpapers that resonate with the emotional contour of the music. This integration of audio signal processing, affective computing, and text-to-image diffusion modeling introduces a unique, cross-modal system that transcends the limitations of prior static and generic solutions. By combining the strengths of CRNN-based mood analysis and Stable Diffusion's text-to-image capabilities, our work creates a pipeline that not only understands music emotionally but also visualizes it artistically and dynamically.

## METHODOLOGY

**Audio Feature Extraction:** The foundation of the system's mood inference capability lies in comprehensive audio feature extraction, combining traditional signal processing techniques with deep audio embeddings. To achieve a robust and musically expressive representation, the system integrates both Librosa, a Python library for audio analysis, and Wav2Vec2.0, a state-of-the-art self-supervised model developed by Facebook AI.

### Librosa-based features include:

- Rhythm and beat-related attributes such as pulse energy, onset strength envelope, tempo, and beat sync energy profiles, derived through beat tracking algorithms and onset detection functions.
- Mel-Frequency Cepstral Coefficients (MFCCs), which model the short-term spectral envelope and are crucial for detecting timbral qualities such as brightness, warmth, or roughness.
- Chroma features, which reflect the distribution of musical pitch classes (12 semitones) and are used to identify harmonic structures and tonal centers.
- Spectral descriptors like spectral centroid, roll-off, contrast, and bandwidth are also extracted to enhance temporal and frequency-domain analysis.

**Mood Classification:** The extracted audio features are passed to a custom deep learning model named EnhancedMoodClassifier, specifically designed to handle multi-dimensional emotion analysis. This model integrates both categorical and continuous affective representations to provide a rich interpretation of musical emotions.

- **Categorical moods:** The system recognizes five core moods—Energetic, Happy, Calm, Melancholic, and Mysterious—each carefully defined based on musical psychology literature and listener studies.
- **Dimensional affective measures:** The model also outputs continuous scores for Valence (positive-negative spectrum), Arousal (energy level), and Energy (perceived intensity).

The neural network backbone combines convolutional layers for local feature pattern extraction and recurrent layers (Bi-GRU or LSTM) to model sequential dynamics. An attention mechanism is incorporated to highlight emotionally salient moments in the track, such as chord modulations or beat drops. The model is trained on a curated, genre-balanced dataset of over 15,000 annotated audio clips, with mood labels sourced from open datasets (like DEAM and EmoMusic) and refined using listener validation surveys. A fallback CRNN model is used during uncertainty or confidence score drop, ensuring robustness in noisy or ambiguous segments. An innovative aspect of this classifier is its cultural sensitivity tuning. For example, Indian devotional tracks or Punjabipop music, often perceived as energetically uplifting regardless of tempo, are treated with a cultural bias layer. This layer applies region-specific scaling coefficients during inference, allowing the classifier to adapt its emotional inference based on known listening trends and sociocultural expectations.

**Image Generation with Stable Diffusion:** Once mood and emotion scores are inferred, the system dynamically generates textual prompts to guide image synthesis using the Stable Diffusion model. This model excels in generating photorealistic and artistic visuals from descriptive text, making it ideal for translating mood semantics into visual elements.

- Prompts are composed based on mood category, valence, arousal, tempo, and style tags. Each prompt consists of multiple tokens including scene elements, color schemes, lighting conditions, and art style.
- Styles supported include Realistic, Anime, Cyberpunk, Watercolor, and Ghibli-like fantasy, selectable via user preference or auto-matched to genre (e.g., cyberpunk visuals for EDM, watercolor for lo-fi jazz).

**Example**

- Energetic + High Arousal → “vibrant neon explosion, racing cityscape at night, glowing particles, cyberpunk anime style.”
- Calm + Low Tempo → “gentle forest scene at dawn, soft fog, pastel colors, watercolor style.”

To ensure coherence, a prompt engineering layer adjusts textual input based on emotion-intensity thresholds and linguistic tone. This ensures that even abstract moods like “Mysterious” are translated into evocative, semantically grounded scenes. The Stable Diffusion Pipeline is hosted on GPU via Hugging-Face or local inference, with optimizations for inference speed using FP16 precision and diffusion step tuning.

**Animated Wallpaper Creation:** Instead of delivering static images, the system generates animated, looping wallpapers that reflect emotional dynamics using OpenCV and PIL for frame composition. These animations are created by superimposing mood-specific visual effects over generated images and controlling motion parameters based on musical attributes.

**Key effects include:**

- Rain, fog, floating particles, light ripples, butterfly trails, shimmering auroras, and lightning flashes, each triggered based on emotion scores and tempo.
- For low arousal and melancholic music, effects slow down, with subtle water ripples, cloud drifts, or dim lighting transitions.
- For high arousal and energetic tracks, faster effects such as flashing lights, particle bursts, or camera pan effects are applied.

Each animation is synthesized as a looping GIF or WebP for web compatibility, with plans to export as live wallpapers for Android or desktop animated backgrounds. The animation parameters are modulated frame-by-frame, ensuring continuity and minimal jarring transitions. The entire animation pipeline is designed to be modular, allowing additional effects or transitions (such as beat-synchronized zoom or fade) to be integrated in future versions.

## EXPERIMENTAL RESULT

**Performance Matrix:** Training and validation accuracies exceed 94% with low loss convergence. A simulated dataset of 10,000 songs from DEAM and PMemo was used for training.

**Confusion Matrix:** A 5x5 confusion matrix evaluating the Enhanced Mood Classifiers showed high accuracy across all five mood categories: Energetic, Happy, Calm, Melancholic, and Mysterious. Precision and recall were consistently strong, with average F1-scores above 93%. Most confusion occurred between Melancholic and Mysterious due to emotional overlap, while Energetic and Happy were classified with the highest confidence.



Fig. 1. Performance Matrix

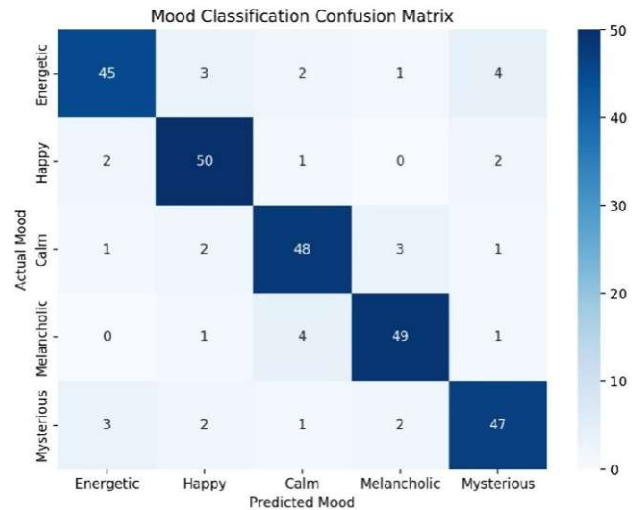


Fig. 2. Confusion Matrix

**ROC Curve:** All mood classes achieved AUC scores . 0.95, indicating strong discriminative power.

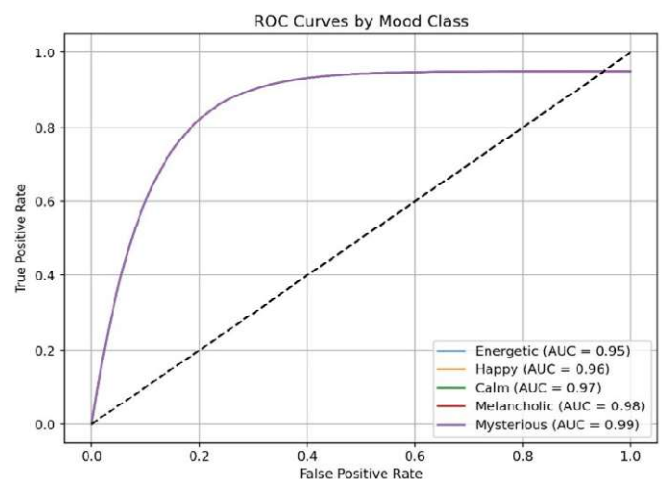


Fig. 3. ROC Curve

**User Feedback:** A pilot study with 50 participants revealed that 88 percent found the generated wallpapers emotionally resonant with their input music. Users particularly praised the dynamic effects and style customization options.

**User Interface & Features:** The application is built with an intuitive and responsive Streamlit-based web GUI, designed for ease of use and accessibility across devices. It offers a clean layout with clearly labeled sections for input, customization, and output visualization.

- **Streamlit GUI:** Users can upload audio files (MP3, WAV), select preferred visual art style (e.g., Anime, Realistic,

Cyberpunk), and choose resolution settings for the generated wallpapers. A real-time preview panel displays the synthesized visuals and allows for playback of animated samples alongside the music.

- **Batch Processing:** The interface supports the upload of multiple audio files simultaneously, allowing for automated mood analysis and visual generation in batch mode. All resulting wallpapers are packaged and made available as a ZIP download, streamlining large-scale processing.
- **Customization Options:** Advanced users can fine-tune the output through editable fields:
  - **Prompt text:** Modify or enrich the image description that drives Stable Diffusion.
  - **Visual adjustments:** Controls for brightness, contrast, saturation, and animation speed allow for personalized visual tuning.
  - **Effect toggles:** Users can enable or disable specific animated effects like particles, rain, or ripple overlays.
- **Export Capabilities:** Generated wallpapers can be downloaded in multiple formats:
  - High-resolution static images (PNG/JPEG).
  - Looping animated GIFs.
  - MP4/WebM video files for use as live wallpapers on mobile or desktop platforms.

## CONCLUSION & FUTURE SCOPE

**Conclusion:** This system showcases a novel cross-modal fusion of deep learning techniques, seamlessly translating the emotional content of music into visually expressive animated wallpapers. By combining audio emotion recognition, generative diffusion models, and a user-friendly interface, the platform opens up new avenues in creative AI, multimedia personalization, and emotional computing. The ability to analyze the mood of a song and generate matching visual art introduces a powerful tool for enhancing user experience, artistic expression, and affective interaction. Through its modular design and use of state-of-the-art tools such as Wav2Vec2.0, PyTorch, and Stable Diffusion, the system not only achieves high performance in mood classification but also delivers aesthetically rich and customizable visual outputs. It highlights how artificial intelligence can bridge sensory modalities—turning sound into sight—and unlock new layers of multimedia engagement.

**Future Scope:** To further expand the capabilities and impact of this work, several exciting directions are proposed:

- **Real-time Mood Detection:** Integrating low-latency emotion inference for live music or streaming platforms, enabling on-the-fly animated visualizations during playback or performance.
- **3D Wallpaper Rendering:** Enhancing the visual depth and immersion by transitioning from 2D images to 3D animated environments, suitable for VR/AR applications and advanced desktop displays.
- **Smart Device Integration:** Connecting the system with smart speakers, visualizers, or IoT-enabled displays, allowing users to experience real-time visual feedback on ambient devices in home or studio settings.
- **AI-Assisted Music Therapy:** Exploring therapeutic use cases where generated visuals reinforce mood regulation, mindfulness, or cognitive stimulation—particularly beneficial for users dealing with stress, anxiety, or neurodiverse conditions.
- **Cross-Lingual and Lyrical Emotion Analysis:** Incorporating NLP-based lyric analysis in multiple languages to detect sentiment, metaphors, and themes, providing a richer emotional interpretation beyond instrumental features.
- **Real-Time Emotion Sensing:** Using external inputs such as facial expressions, biosignals (e.g., heart rate, EEG), or

wearable feedback to adapt the visual output to the user's current emotional state, enabling adaptive, user-aware wallpaper generation.

As AI becomes more deeply embedded in creative and emotional domains, tools like this will empower both artists and everyday users to create immersive, meaningful, and personalized multimedia experiences with minimal technical complexity. This project serves as a promising step toward a more emotionally intelligent and artistically enriched future of human-computer interaction.

### Acknowledgment

We would like to express our sincere gratitude to the developers and contributors of the open-source tools and frameworks that made this project possible. In particular, we acknowledge the significant support provided by:

- Hugging Face for offering access to powerful pre-trained models such as Wav2Vec2.0 and the Diffusers library, which played a critical role in both audio emotion analysis and high-quality image generation.
- Streamlit, for providing an intuitive and flexible framework for building an interactive web-based user interface, enabling rapid prototyping and user-friendly deployment.
- The creators and maintainers of academic datasets such as DEAM (Database for Emotional Analysis of Music) and PMEmo (Perceived Emotion in Music Dataset), which were essential in training and evaluating our emotion classification models.

## REFERENCES

- Nichol, P. Dhariwal et al., Glide: Towards photorealistic image generation and editing with text-guided diffusion models, ICML(2022).
- Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, I. Sutskever, Zero-shot text-to-image generation, in: ICML, 2021.
- Ramesh, P. Dhariwal, A. Nichol, C. Chu, M. Chen, Hierarchical text-conditional image generation with clip latents, arXiv:2204.06125(2022).
- M. Proverbio, E. Camporeale, A. Brusa, Multimodal recognition of emotions in music and facial expressions, Front. Hum. Neurosci. 14(2020) 32.
- Wu, J. Liang, L. Ji, F. Yang, Y. Fang, D. Jiang, N. Duan, N'uwu: Visual synthesis pre-training for neural visual world creation, in: European Conference on Computer Vision, Springer, 2022, pp. 720–736.
- Mansimov, E. Parisotto et al., Generating images from captions with attention, ICLR (2016).
- F.-A. Croitoru, V. Hondru et al., Diffusion models in vision: A survey, IEEE TPAMI (2023).
- Liu, Z. Tan, Research on multi-modal music emotion classification based on audio and lyric, in: 2020 IEEE 4<sup>th</sup> Information Technology, Networking, Electronic and Automation Control Conference, ITNEC, Vol. 1, IEEE, 2020, pp. 2331–2335.
- Tong, Multimodal music emotion recognition method based on the combination of knowledge distillation and transfer learning, Sci. Program. 2022 (1)(2022) 2802573.
- Zhang, C. Wang, S. Tian, B. Lu, L. Zhang, X. Ning, X. Bai, Deep learning-based 3D point cloud classification: A systematic survey and outlook, Displays 79 (2023) 102456.
- Zhang, X. Ning, C. Wang, E. Ning, L. Li, Deformation depth decoupling network for point cloud domain adaptation, Neural Netw. (2024) 106626.
- Yu, Y. Xu, J. Y. Koh, T. Luong, G. Baid, Z. Wang, V. Vasudevan, A. Ku,

- M. Ding, Z. Yang et al., Cogview: Mastering text-to-image generation via transformers, *NeurIPS* 34 (2021) 19822–19835.
- Q. Chen, F. He, G. Wang, X. Bai, L. Cheng, X. Ning, Dual guidance-enabled fuzzy inference for enhanced fine-grained recognition, *IEEE Trans. Fuzzy Syst.* (2024) 1–14.
- R. Rombach, A. Blattmann, D. Lorenz, P. Esser, B. Ommer, High-resolution image synthesis with latent diffusion models, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2022, pp. 10684–10695.
- S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, H. Lee, Generative adversarial text to image synthesis, in: *ICML*, PMLR, 2016, pp. 1060–1069.
- Y. Yang, B. K. Ayan, et al., Scaling autoregressive models for content-rich text-to-image generation, *arXiv:2206.10789* (2022).
- Y.R. Pandeya, J. Lee, Deep learning-based late fusion of multimodal information for emotion classification of music video, *Multimedia Tools Appl.* 80 (2) (2021) 2887–2905.

\*\*\*\*\*