

---

# V2Meow: Meowing to the Visual Beat via Music Generation

---

Anonymous Author(s)

Affiliation

Address

email

## 1 Introduction

The majority of prior research in video-to-music generation has concentrated on designing complex rhythmic feature extractors to model the physical correspondence between video content and music. For instance, generating music synchronized with dance movements or reconstructing instrumental music based on changes in human poses in silent instrument performance videos (14; 15; 13; 11; 6). Consequently, these approaches are typically tailored to specific visual scenarios, and cannot be generalized to arbitrary video input types, e.g., vlogs or slideshows of images. In contrast, our study delves into the challenge of generating contextually relevant and high-quality background music for a broad spectrum of video input types. Importantly, we achieve this by conditioning the music generation solely on the visual information provided by video frames, without explicitly modeling domain-specific rhythmic or semantic relationships. We further hypothesize that with sufficient data and scale, the generation model is capable of learning the intrinsic video-music correspondence directly from easily accessible music videos and generating relevant background music for unseen video content types through zero-shot transfer.

We propose a video-to-music generation system called V2Meow that can generate high-quality music audio for a diverse range of video input types based on a multi-stage autoregressive model, without the need to explicitly model the rhythmic or semantic video-music correspondence. Compared to previous video to music generation work, the video and text prompts are incorporated as a single stream of embedding inputs and fed into the Transformer with feature-specific adaptors. Trained on O(100K) music audio clips paired with video frames mined from in-the-wild music videos, V2Meow is competitive with previous domain-specific models when evaluated in a zero-shot manner. V2Meow can synthesize high-fidelity music audio waveform solely by conditioning on pre-trained general-purpose visual features extracted from video frames, with optional style control via text prompts. Through both qualitative and quantitative evaluations, we verify that our model outperforms various existing music generation systems in terms of visual-audio correspondence and audio quality.

## 2 Method

Inspired by MusicLM (2), we take a multi-stage autoregressive language modeling approach (Figure 1) to condition music generation on video frames with optional high-level control over the style of the generated music through text prompts.

**Feature Representations.** For audio representations, we adopt the SoundStream tokenizer for acoustic tokens modeling and w2v-BERT tokenizer for semantic tokens modeling, both of which are pre-trained on the Free Music Archive (FMA) dataset (3). For all visual features, we use frame rate at 1 fps, following the standard on MV100K (1). For visual feature to music semantic tokens modeling, we use encoder-decoder Transformer and use 10-second random crops of the music video for visual to music semantic tokens modeling and semantic tokens to coarse acoustic tokens modeling. During inference, we take 10-second silent video as input and generate 10-second music clip.

**Optional Style Control.** To incorporate the control signal, we simply feed the MuLan audio embedding (8) as an additional input with sequence length be one to the Transformer encoder along with the visual features in the first stage. Both Mulan audio embedding and visual features are

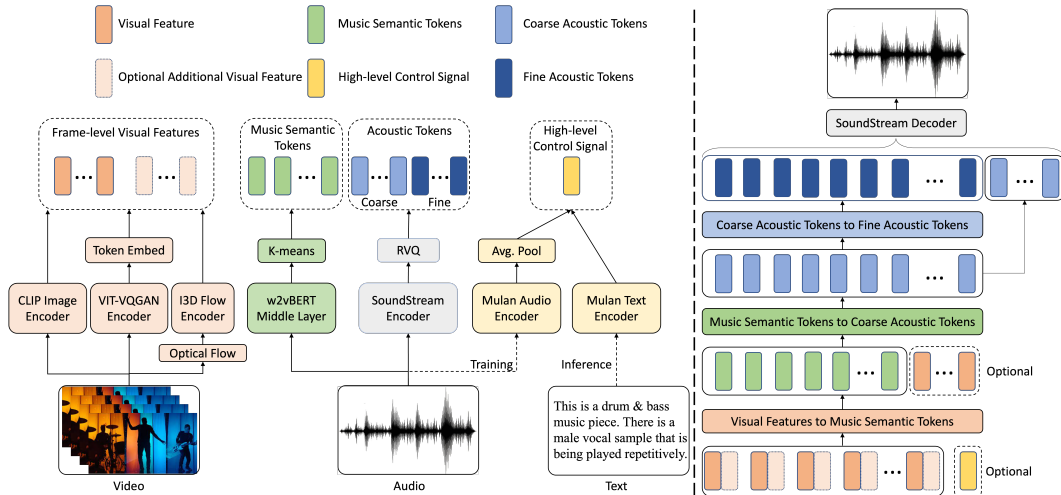


Figure 1: V2Meow architecture overview: (left) Feature extraction pipeline for video, audio and text representations. (right) Overview of multi-stage video to music modeling.

40 projected to the same feature dimension. At inference, we instead use the MuLan text embedding  
 41 with the visual features to generate the semantic tokens.

42 **Datasets.** Following (12), we filtered a public available video dataset (1) to 110k videos with the  
 43 label Music Videos and refer to it as MV100K. The training and validation datasets were split into an  
 44 80:20 ratio. We trained the Stage 1 model and Stage 2 model on these O(100K) music videos and  
 45 refer to it as MV100K.

### 46 3 Evaluation and Results

47 We conduct a comprehensive evaluation of the music generation methods across three distinct datasets.  
 48 To assess the quality of the generated music, we employ a multifaceted evaluation framework that  
 49 encompasses both quantitative and qualitative metrics. Quantitative measures include evaluations of  
 50 audio quality, rhythm synchronization, and text alignment. Additionally, we employ subjective metrics  
 51 like visual relevance and music preference, which are best assessed through human evaluations.

52 **Video conditional music generation.** Since there is no open-source video to music in audio  
 53 waveform, we compare our V2Meow model against the state-of-the-arts video-driven symbolic  
 54 music representations-based model CMT (4) on the test partition of the MV100K. In terms of  
 55 visual relevance and music quality, V2Meow significantly outperforms CMT by a large margin. For  
 56 MV100K, we observe that adding visual input at the acoustic modeling stage significantly improves  
 57 both audio quality related metrics. We further observe that the combination of CLIP and I3D Flow  
 58 features yield best music generation quality overall.

59 **Video and text conditional music generation.** We compare V2Meow with text-to-music generation  
 60 models like MUBERT (10) and Riffusion (5) on latest MusicCaps dataset (2), a subset of AudioSet (7)  
 61 that contains about 5.5k human annotated text captions, music, and video pairs. With video frames  
 62 as additional control, our approach outperforms Riffusion and MUBERT in visual relevance by  
 63 20-30%. It is worth to note that while our V2Meow model only use video-level MuLan embedding  
 64 and trained on a O(100K) music videos, we still achieve better audio quality and text adherence than  
 65 pure text-to-music generative model.

66 **Dance to music generation.** Evaluation on 20 dance videos in the test split of AIST++ (9) demon-  
 67 strates that V2Meow can achieve comparable performances to domain-specific dance-to-music  
 68 generation baselines (14; 15), as measured by beat coverage and beat hit. The evaluation is in  
 69 zero-shot fashion without any fine-tuning on the AIST++ train split, and only video frames are used  
 70 for modeling while no motion data is involved.

## 71 4 Ethical Implications

72 Controllable generative models such as V2Meow can serve as the foundation for new tools, technolo-  
73 gies, and practices for content creators. While our motivation is to support creators to enrich their  
74 creative pursuits, we acknowledge that large generative models learn to imitate patterns and biases  
75 inherent in the training sets, and in our case, the model can propagate the potential biases built in the  
76 video and music corpora used to train our models.

77 Such biases can be hard to detect as they manifest in often subtle, unpredictable ways, which are not  
78 fully captured by our current evaluation benchmarks. Demeaning or other harmful language may be  
79 generated in model outputs, due to learned associations or by chance. A thorough analysis of our  
80 training dataset shows that the genre distribution is skewed towards a few genres, and within each  
81 genre, gender, age or ethical groups are not represented equally. For example, male is dominant in  
82 hip-hop and heavy metal genre. These concerns extend to learned visual-audio associations, which  
83 may lead to stereotypical associations between video content (i.e. people, body movements/dance  
84 styles, locations, objects) and a narrow set of musical genres; or to demeaning associations between  
85 choreography in video content and audio output (i.e. minstrelsy, parody, miming). ML fairness  
86 testing is required to understand the likelihood of these patterns in any given model and effectively  
87 intervene in them.

## 88 References

- 89 [1] Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan,  
90 and Sudheendra Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. *arXiv*  
91 *preprint arXiv:1609.08675*, 2016.
- 92 [2] Andrea Agostinelli, Timo I. Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing  
93 Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, Matthew Sharifi, Neil Zeghidour, and C. Frank.  
94 Musiclm: Generating music from text. *ArXiv*, abs/2301.11325, 2023.
- 95 [3] Michaël Defferrard, Kirell Benzi, Pierre Vandergheynst, and Xavier Bresson. FMA: A dataset for music  
96 analysis. In *18th International Society for Music Information Retrieval Conference (ISMIR)*, 2017.
- 97 [4] Shangzhe Di, Zeren Jiang, Si Liu, Zhaokai Wang, Leyan Zhu, Zexin He, Hongming Liu, and Shuicheng  
98 Yan. Video background music generation with controllable music transformer. In *Proc. of the ACM*  
99 *International Conference on Multimedia*, 2021.
- 100 [5] Seth\* Forsgren and Hayk\* Martiros. Riffusion - Stable diffusion for real-time music generation. 2022.
- 101 [6] Chuang Gan, Deng Huang, Peihao Chen, Joshua B. Tenenbaum, and Antonio Torralba. Foley music:  
102 Learning to generate music from videos. In *Proc. of the European Conference on Computer Vision (ECCV)*,  
103 2020.
- 104 [7] Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore,  
105 Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In  
106 *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 776–780.  
107 IEEE, 2017.
- 108 [8] Qingqing Huang, Aren Jansen, Joonseok Lee, Ravi Ganti, Judith Yue Li, and Daniel PW Ellis. MuLan:  
109 A joint embedding of music audio and natural language. In *Proc. of the International Society for Music*  
110 *Information Retrieval Conference (ISMIR)*, 2022.
- 111 [9] Ruilong Li, Shan Yang, David A. Ross, and Angjoo Kanazawa. Ai choreographer: Music conditioned 3d  
112 dance generation with aist++, 2021.
- 113 [10] Mubert-Inc. Mubert. <https://mubert.com/>, <https://github.com/mubertai/mubert-text-to-music>. 2022.
- 114 [11] Kun Su, Xiulong Liu, and Eli Shlizerman. Audeo: Audio generation for a silent performance video.  
115 *Advances in Neural Information Processing Systems*, 33, 2020.
- 116 [12] Dídac Surís, Carl Vondrick, Bryan Russell, and Justin Salamon. It’s time for artistic correspondence in  
117 music and video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
118 pages 10564–10574, 2022.
- 119 [13] Jiashuo Yu, Yaohui Wang, Xinyuan Chen, Xiao Sun, and Y. Qiao. Long-term rhythmic video soundtracker.  
120 *ArXiv*, abs/2305.01319, 2023.
- 121 [14] Ye Zhu, Kyle Olszewski, Yuehua Wu, Panos Achlioptas, Menglei Chai, Yan Yan, and S. Tulyakov.  
122 Quantized gan for complex music generation from dance videos. *ArXiv*, abs/2204.00604, 2022.
- 123 [15] Ye Zhu, Yuehua Wu, Kyle Olszewski, Jian Ren, S. Tulyakov, and Yan Yan. Discrete contrastive diffusion  
124 for cross-modal and conditional generation. *ArXiv*, abs/2206.07771, 2022.