

Cross-Modal Learning for Music-to-Music-Video Description Generation

Zhuoyuan Mao¹ Mengjie Zhao¹ Qiyu Wu¹ Zhi Zhong¹
Wei-Hsiang Liao² Hiromi Wakaki¹ Yuki Mitsufuji^{1,2}

¹Sony Group Corporation ²Sony AI

{zhuoyuan.mao, mengjie.zhao, qiyu.wu, zhi.zhong,
weihsiang.liao, hiromi.wakaki, yuhki.mitsufuji}@sony.com

Abstract

Music-to-music-video generation is a challenging task due to the intrinsic differences between the music and video modalities. The advent of powerful text-to-video diffusion models has opened a promising pathway for music-video (MV) generation by first addressing the music-to-MV description task and subsequently leveraging these models for video generation. In this study, we focus on the MV description generation task and propose a comprehensive pipeline encompassing training data construction and multimodal model fine-tuning. We fine-tune existing pre-trained multimodal models on our newly constructed music-to-MV description dataset based on the Music4All dataset, which integrates both musical and visual information. Our experimental results demonstrate that music representations can be effectively mapped to textual domains, enabling the generation of meaningful MV description directly from music inputs. We also identify key components in the dataset construction pipeline that critically impact the quality of MV description and highlight specific musical attributes that warrant greater focus for improved MV description generation.

1 Introduction

Generating a music-video (MV) to match a given piece of music is a challenging task due to the inherent differences between the music and video modalities. Despite the challenges, MV generation holds significant potential for enhancing the music experience by providing a visual narrative that aligns with the music’s tone, style, and mood, offering a more immersive and engaging way for audiences to connect with the music. Compared to generating music or audio from a given video (Tian et al., 2024; Kang et al., 2024), the reverse task is more complex, as the video modality typically conveys richer spatial and temporal information than music. However, with the advent of text-to-

video diffusion models (Yang et al., 2024; Polyak et al., 2024; Kong et al., 2024), videos can now be generated from textual descriptions. This development enables MV generation to be divided into two subtasks: (1) music-to-MV description generation and (2) MV description-to-MV generation. As illustrated in Fig. 1, MV descriptions can be further refined using large language models (LLMs) like GPT (OpenAI, 2023) to fit specific text-to-video models (Khachatryan et al., 2023). *In this study, we focus on the first task: generating MV descriptions from music.*

To this end, we propose a practical pipeline for data construction and model training to generate meaningful MV descriptions based on music inputs. Additionally, we explore methods to enhance the alignment of the generated descriptions to the given music. Specifically, we investigate the impact of various data sources—such as music, music genre tags, MV type tags, and lyrics understanding text—on the quality of the generated MV descriptions when fine-tuning multimodal LLMs like NExT-GPT (Wu et al., 2024). As shown in Fig. 1, our approach first leverages existing music understanding models (Zhao et al., 2024; Mao et al., 2025) to extract lyrics understanding text. We then fine-tune a multimodal LLM to process these diverse inputs and generate MV descriptions. The training data is constructed from gold-standard MVs, incorporating music-related information to enhance the connection between music and the generated descriptions. Unlike prior studies on MV generation, such as ViPE (Shahmohammadi et al., 2023), which focused solely on lyrics as input, our work emphasizes leveraging multiple modalities and evaluates the effectiveness of various combinations of input data in connecting multimodal representations for MV description generation.

To facilitate this study, we construct a music-to-MV description training and evaluation dataset using the Music4All dataset (Santana et al., 2020).

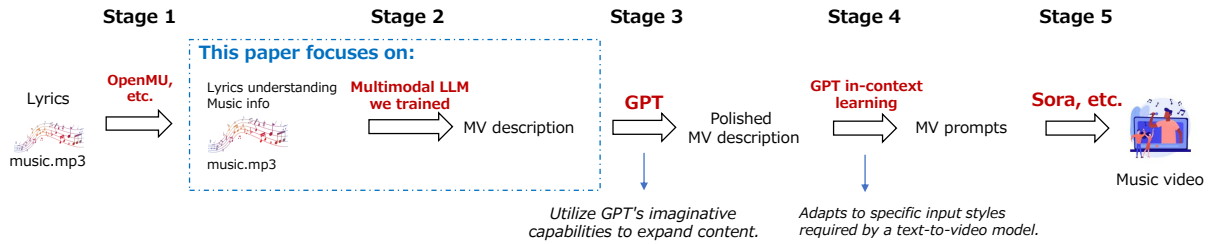


Figure 1: **Pipeline of music-to-MV generation.** We focus on multimodal model training of Stage 2 in this study.

Empirical results on the NExT-GPT baseline and multimodal LLMs fine-tuned with our dataset demonstrate that meaningful MV descriptions can be generated from music and music-related textual inputs after multimodal fine-tuning. An ablation study on different combinations of input sources, including music, music genre tags, MV type tags, and lyrics understanding text, reveals that music and MV type tags are key components for high-quality MV description generation. While music genre tags and lyrics understanding text also contribute positively, they can be used interchangeably. Our findings can contribute to future study on enhancing MV descriptions and temporal alignment between music, lyrics, and the generated MV.

2 Proposed Method

In this section, we present the pipeline proposed for training a multimodal LLM specifically tailored to the music-to-MV description generation task. For the first time, our pipeline incorporates a broader range of musical information beyond lyrics as inputs, aiming to enrich the holistic understanding of the music. Additionally, we introduce strategies to ensure the generated MV descriptions are more closely aligned with the musical inputs. The curated dataset is then utilized to fine-tune a multimodal LLM for performing the MV description generation task.

2.1 Data Construction

This section outlines our proposed pipeline for constructing training and evaluation datasets for the music-to-MV description generation task.

2.1.1 MV Datasets

We construct our datasets based on the Music4All dataset (Santana et al., 2020), which comprises approximately 100k music clips paired with corresponding MVs and enriched with metadata such as energy, valence, and genre. To enhance the dataset, we leverage the OpenMU model (Zhao et al., 2024)

to generate lyrics understanding text for all music clips in Music4All. This process effectively interprets the lyrics for each piece of music, providing concise textual information related to the lyrics. Additionally, we filter out MVs that consist solely of static images rather than actual video footage. After filtering, the final dataset includes 56,446 samples, 55,000 for training and 1,446 for testing.

2.1.2 Construction of Input Data for Music and Associated Information

After preparing the training and evaluation splits of music clips, MVs, lyrics understanding text, and metadata from the Music4All dataset, we curate various data types as inputs for the MV description generation task. To incorporate richer musical information across different modalities, we include music genre tags and lyrics understanding text as inputs in addition to the music clips. Moreover, to refine the output MV descriptions and make the task less open-ended, we specify the style of the output by providing MV type tags. These tags are assigned to the MV clips using GPT-4o mini (OpenAI, 2023) and include ten category candidates: Live Performance, Lyric Video, Animation, Story Narrative, Artistic/Abstract, Dance Performance, Behind-the-Scenes, Nature/Scenic, Static/Dynamic Picture Montage, and Cinematic Drama.¹

As shown in Fig. 2, the four types of inputs are used to train the multimodal LLM, guided by a fixed instruction: “Generate a concise video prompt that captures the essence of the MV, incorporating the music’s tone, style, and lyrical themes. The prompt should reflect the specified MV type and align with the music genre to ensure stylistic coherence for guiding a text-to-video model.”

2.1.3 Construction of Output Data for MV Description

The output MV descriptions should provide rich content detailing the visual elements of the MV

¹Generated based on suggestions from GPT-4o mini.

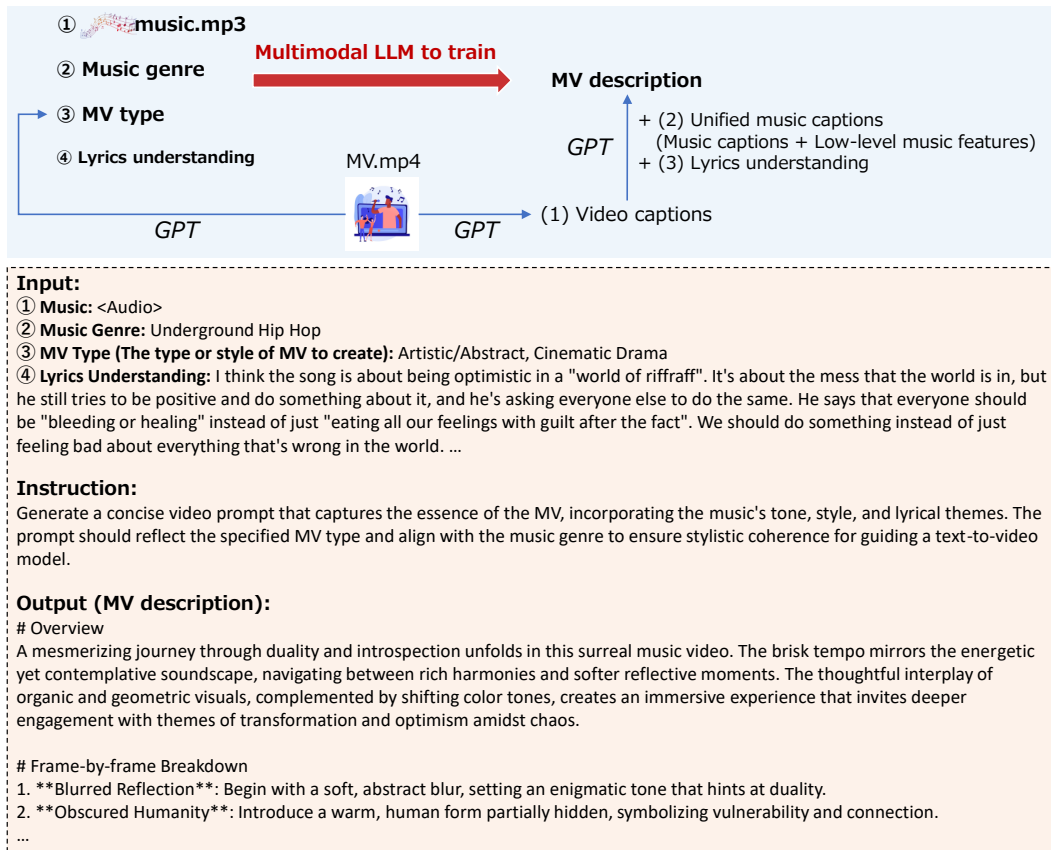


Figure 2: Process for creating music-to-MV description training datasets (top) and an example of utilizing the generated data to train music-to-text LLMs (bottom).

while remaining closely tied to musical features, such as tempo, downbeats, and high-level characteristics like the mood conveyed by the music. To achieve this, we first utilize GPT-4o mini to caption MV clips and extract relevant visual contexts. Next, we refine these captions using GPT-4o mini again, integrating key musical features, including music captions, low-level musical attributes, and lyrics understanding. Music captions and lyrics understanding texts are generated using the OpenMU music understanding model, while low-level musical features are extracted with open-source tools (Böck et al., 2016), following the methodology of LLark (Gardner et al., 2024). The constructed MV description dataset includes two main components: an overview and a frame-by-frame breakdown, with frame captions extracted at two-second intervals for each 30-second MV clip. Examples of music captions, low-level music features, and a complete version of an MV description are provided in Appendix A.

2.2 Multimodal Model Training

We utilize NExT-GPT (Wu et al., 2024), an any-to-any multimodal training framework, to fine-tune our model using the constructed music-to-MV description datasets. Following NExT-GPT’s methodology, the fine-tuning process is divided into multiple stages. In the first stage, we fine-tune only the adaptor between the ImageBind (Girdhar et al., 2023) encoder and the Vicuna LLM (Zheng et al., 2023) utilizing the music captioning task. In the second stage, we simultaneously fine-tune the adaptor and apply LoRA (Hu et al., 2022) fine-tuning to Vicuna with the constructed music-to-MV description dataset. As illustrated in Fig. 2, the input data including the music clip is sequentially formatted, followed by a fixed instruction. The model is trained to generate MV descriptions comprising an overall summary and frame-by-frame breakdowns. We fine-tune for 5 and 2 epochs in the first and second stages, respectively, utilizing a learning rate of $1e-4$ and a batch size of 2. Training is conducted on 2 NVIDIA A6000 GPUs. For LoRA, the rank and alpha are both set to 32, following NExT-GPT.

Model	BLEU-1	BLEU	ROUGE-P	ROUGE-R	ROUGE-F1	BERT-P	BERT-R	BERT-F1
<i>Baseline</i>								
NExT-GPT (Wu et al., 2024)	8.3	0.2	20.7	9.2	11.8	80.9	76.5	78.6
<i>Main results</i>								
①+②+③+④	42.9	14.6	22.9	23.2	22.7	87.4	86.4	86.9
<i>Ablation study</i>								
②+③+④	42.5	14.4	22.4	22.8	22.3	87.2	86.2	86.7
①+②+③	43.6	14.7	23.0	23.5	22.9	87.4	86.7	87.0
①+③+④	42.8	14.5	22.8	23.2	22.7	87.3	86.4	86.9
①+②+④	42.2	14.1	21.7	22.5	21.8	86.9	86.2	86.5
②+③	41.8	14.0	21.8	22.4	21.8	87.2	86.1	86.6
①+③	42.4	14.4	22.3	22.8	22.2	87.2	86.2	86.6
①+④	41.3	13.8	21.4	22.4	21.6	86.8	86.0	86.4
<i>Sanity check (w/o inputs, solely w/ instructions during inference)</i>								
①+②+③+④	39.3	13.2	20.2	22.5	21.0	85.8	85.6	85.7
①+④	39.7	12.5	20.3	20.8	20.3	86.1	85.6	85.9

Table 1: Results of MV description generation on the Music4All dataset. We provide BLEU-1 and BLEU-4 scores for BLEU, along with precision, recall, and F1 scores for both ROUGE-L and BERT-score. ①, ②, ③, and ④ represent music, music genre tags, MV type tags, and lyrics understanding text, respectively. The top three values in each metric are highlighted in **bold**.

3 Evaluation

Using the 1,446 test samples from our constructed dataset, we evaluate the generated MV descriptions with BLEU (Papineni et al., 2002), ROUGE-L (Lin, 2004), and BERT-score (Zhang et al., 2020), considering different combinations of inputs: ① music, ② music genre tags, ③ MV type tags, and ④ lyrics understanding text. Additionally, we present several MV frames generated by Text2Video-Zero (Khachatryan et al., 2023) to test the feasibility of the entire proposed MV generation pipeline in Appendix B, using the ground-truth MV descriptions annotated by us as input.

3.1 Main Results

As shown in Table 1, our proposed pipeline for music-to-MV description generation achieves significant improvements over the NExT-GPT baseline after fine-tuning for a specific music domain. This demonstrates that, with the proposed datasets and pipeline, music can be effectively mapped to the text domain. Comparing the main results with sanity checks that remove all inputs during inference (leaving only a fixed instruction), we observe that our carefully designed inputs for music-related information substantially contribute to the quality of the generated MV descriptions. Interestingly, after training, the model can generate reasonable MV descriptions even without any inputs, suggesting that the NExT-GPT model successfully adapts to the MV description generation downstream task.

3.2 Ablation Study

Through an ablation study exploring different combinations of data sources, we find that settings ①+②+③ and ①+③+④ achieve comparable or even slightly better performance to the full data combination (①+②+③+④). This suggests that the contributions of music genre tags (②) and lyrics understanding text (④) are interchangeable, without providing additional benefits when used together. Observing the results of ①+③, we note that music genre tags (②) and lyrics understanding (④) positively impact the results and are not redundant inputs. When comparing the top three performing settings (①+②+③, ①+③+④, and ①+②+③+④) with the combinations ②+③+④ and ①+②+④, we observe a significant performance drop. This highlights the importance of including both music (①) and MV type tags (③). Seeing the results of ①+④, the simultaneous inclusion of music genre tags (②) and MV type tags (③) yields consistent improvement across all metrics. Moreover, the results of ②+③ demonstrate that even with simple tags for music and MV, the model can generate reasonable MV descriptions, suggesting future opportunities to enhance the model’s performance by leveraging finer-grained features such as temporal alignment between lyrics and musical waves.

4 Conclusion

In this study, we explored data construction and multimodal training pipelines for the music-to-MV

description task, with the goal of building robust base models for the broader music-to-MV generation task. Our results on the constructed Music4All dataset highlight key data sources that significantly impact the quality of MV descriptions. Future work could extend our proposed dataset construction pipeline to additional music domains. Additionally, exploring more effective methods to align MV descriptions or MVs with the corresponding music could pave the way for improved performance in this challenging task.

Limitations

The proposed approach has several limitations: (1) The pipeline was evaluated on a single constructed dataset. Testing on additional datasets could strengthen the claims made in this paper. (2) The pipeline focuses on converting music into MV descriptions for MV generation tasks, but relying solely on text descriptions may overlook important information necessary for effective MV generation. (3) Inputs were limited to music, music genre tags, MV type tags, and lyrics understanding text, while other features that could significantly enhance MV descriptions may not have been considered. (4) The data construction pipeline depends on LLMs for captioning, and the choice of LLMs could influence the quality of the generated MV descriptions.

Ethical Considerations

In this study, we utilized only publicly available datasets and models to fine-tune the music-to-MV description generation task, ensuring no copyright issues. While our experiments focused on MV description generation, it is important to acknowledge that the fine-tuned models may produce potentially risky hallucinations. Users should use the generated content with caution, understand the risks associated with LLM-generated outputs, and implement content safety checks as post-processing. Although debiasing fine-tuning could help address these issues, it falls outside the scope of this work. Additionally, caution is needed when using text-to-video models based on the generated MV descriptions, ensuring that no illegal content, such as unauthorized human identities or privacy violations, is included.

References

Sebastian Böck, Filip Korzeniowski, Jan Schlüter, Florian Krebs, and Gerhard Widmer. 2016. *madmom*:

A new python audio and music signal processing library. In *Proceedings of the 2016 ACM Conference on Multimedia Conference, MM 2016, Amsterdam, The Netherlands, October 15-19, 2016*, pages 1174–1178. ACM.

Joshua Patrick Gardner, Simon Durand, Daniel Stoller, and Rachel M. Bittner. 2024. *Lark: A multi-modal instruction-following language model for music*. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.

Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Manohar Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. 2023. *Imagebind one embedding space to bind them all*. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 15180–15190. IEEE.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. *Lora: Low-rank adaptation of large language models*. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.

Jaeyong Kang, Soujanya Poria, and Dorien Herremans. 2024. *Video2music: Suitable music generation from videos using an affective multimodal transformer model*. *Expert Syst. Appl.*, 249:123640.

Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. 2023. *Text2video-zero: Text-to-image diffusion models are zero-shot video generators*. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pages 15908–15918. IEEE.

Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, Kathrina Wu, Qin Lin, Junkun Yuan, Yanxin Long, Aladdin Wang, Andong Wang, Changlin Li, Duojun Huang, Fang Yang, Hao Tan, Hongmei Wang, Jacob Song, Jiawang Bai, Jianbing Wu, Jinbao Xue, Joey Wang, Kai Wang, Mengyang Liu, Pengyu Li, Shuai Li, Weiyan Wang, Wenqing Yu, Xincheng Deng, Yang Li, Yi Chen, Yutao Cui, Yuanbo Peng, Zhentao Yu, Zhiyu He, Zhiyong Xu, Zixiang Zhou, Zunnan Xu, Yangyu Tao, Qinglin Lu, Songtao Liu, Daquan Zhou, Hongfa Wang, Yong Yang, Di Wang, Yuhong Liu, Jie Jiang, and Caesar Zhong. 2024. *Hunyuanvideo: A systematic framework for large video generative models*. *CoRR*, abs/2412.03603.

Chin-Yew Lin. 2004. *ROUGE: A package for automatic evaluation of summaries*. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Zhuoyuan Mao, Mengjie Zhao, Qiyu Wu, Hiromi Wakaki, and Yuki Mitsufuji. 2025. *Deepresonance*:

- Enhancing multimodal music understanding via music-centric multi-way instruction tuning. *CoRR*, abs/2502.12623.
- OpenAI. 2023. *GPT-4 technical report*. *CoRR*, abs/2303.08774.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. *Bleu: a method for automatic evaluation of machine translation*. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra, Animesh Sinha, Ann Lee, Apoorv Vyas, Bowen Shi, Chih-Yao Ma, Ching-Yao Chuang, David Yan, Dhruv Choudhary, DingKang Wang, Geet Sethi, Guan Pang, Haoyu Ma, Ishan Misra, Ji Hou, Jialiang Wang, Kiran Jagadeesh, Kunpeng Li, Luxin Zhang, Mannat Singh, Mary Williamson, Matt Le, Matthew Yu, Mitesh Kumar Singh, Peizhao Zhang, Peter Vajda, Quentin Duval, Rohit Girdhar, Roshan Sumbaly, Sai Saketh Rambhatla, Sam S. Tsai, Samaneh Azadi, Samyak Datta, Sanyuan Chen, Sean Bell, Sharadh Ramaswamy, Shelly Sheynin, Siddharth Bhattacharya, Simran Motwani, Tao Xu, Tianhe Li, Tingbo Hou, Wei-Ning Hsu, Xi Yin, Xiaoliang Dai, Yaniv Taigman, Yaqiao Luo, Yen-Cheng Liu, Yi-Chiao Wu, Yue Zhao, Yuval Kirstain, Zecheng He, Zijian He, Albert Pumarola, Ali K. Thabet, Artsiom Sanakoyeu, Arun Mallya, Baishan Guo, Boris Araya, Breena Kerr, Carleigh Wood, Ce Liu, Cen Peng, Dmitry Vengertsev, Edgar Schönfeld, Elliot Blanchard, Felix Juefei-Xu, Fraylie Nord, Jeff Liang, John Hoffman, Jonas Kohler, Kaolin Fire, Karthik Sivakumar, Lawrence Chen, Licheng Yu, Luya Gao, Markos Georgopoulos, Rashel Moritz, Sara K. Sampson, Shikai Li, Simone Parmeggiani, Steve Fine, Tara Fowler, Vladan Petrovic, and Yuming Du. 2024. *Movie gen: A cast of media foundation models*. *CoRR*, abs/2410.13720.
- Igor André Pegoraro Santana, Fabio Pinhelli, Juliano Donini, Leonardo Gabiato Catharin, Rafael Biazus Mangolin, Yandre Maldonado e Gomes da Costa, Valéria Delisandra Feltrim, and Marcos Aurélio Domingues. 2020. *Music4all: A new music database and its applications*. In *2020 International Conference on Systems, Signals and Image Processing, IWS-SIP 2020, Niterói, Brazil, July 1-3, 2020*, pages 399–404. IEEE.
- Hassan Shahmohammadi, Adhiraj Ghosh, and Hendrik Lensch. 2023. *ViPE: Visualise pretty-much everything*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5477–5494, Singapore. Association for Computational Linguistics.
- Zeyue Tian, Zhaoyang Liu, Ruibin Yuan, Jiahao Pan, Xiaoqiang Huang, Qifeng Liu, Xu Tan, Qifeng Chen, Wei Xue, and Yike Guo. 2024. *Vidmuse: A simple video-to-music generation framework with long-short-term modeling*. *CoRR*, abs/2406.04321.
- Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. 2024. *Next-gpt: Any-to-any multimodal LLM*. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, Da Yin, Xiaotao Gu, Yuxuan Zhang, Weihang Wang, Yean Cheng, Ting Liu, Bin Xu, Yuxiao Dong, and Jie Tang. 2024. *Cogvideox: Text-to-video diffusion models with an expert transformer*. *CoRR*, abs/2408.06072.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. *Bertscore: Evaluating text generation with BERT*. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Mengjie Zhao, Zhi Zhong, Zhuoyuan Mao, Shiqi Yang, Wei-Hsiang Liao, Shusuke Takahashi, Hiromi Wakaki, and Yuki Mitsufuji. 2024. *Openmu: Your swiss army knife for music understanding*. *CoRR*, abs/2410.15573.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. *Judging llm-as-a-judge with mt-bench and chatbot arena*. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

A Examples of How to Construct MV Description

We first extract low-level music features, including tempo, key, downbeats, and chords, using the open-source tool madmom (Böck et al., 2016). Based on these features and the textual captions of each music piece, we employ GPT-4o mini² to generate unified music captions that seamlessly integrate all the musical information into natural, coherent sentences, as illustrated in Figure 3. Subsequently, we prompt GPT-4o mini again to construct MV descriptions by combining the video captions of each gold-standard MV, the unified music captions, and the lyrics understanding text (see Figure 2). The resulting MV descriptions incorporate both visual and musical content, making them better suited for reconstructing the original MV.

B Generating Video Frames using MV Description

Figure 5 showcases frames generated by the Text2Video-Zero (Khachatryan et al., 2023) model using the gold-standard MV description example provided in Figure 4. When compared to the original MV frames, we observe that even with only textual descriptions, the text-to-video model can produce content closely aligned with the intended visuals, such as the abstract geometric shapes in frames #3 to #5 and the mirrored sky in frames #6 and #7. This demonstrates the feasibility of our proposed pipeline for MV generation, as illustrated in Figure 1. However, challenges remain, particularly in accurately generating complex elements like multi-layered imagery and human faces using current text-to-video models. Addressing these limitations will be crucial for future advancements in this domain.

²<https://platform.openai.com/docs/models#gpt-4o-mini>

Input: Music captions, low-level music features

Music captions: The music clip features a very fast tempo with medium energy, creating a sense of urgency in the rhythm. It has a low valence, conveying a somber or melancholic mood. The overall danceability is moderate, suggesting a balance between movement and reflection in the music. The combination of the brisk tempo, medium energy, and low valence gives the clip a unique sense of momentum tempered with emotional depth.

Low-level music features: {'tempo': [[85.71428571428571, 0.31086550309221533], [43.16546762589928, 0.16042116950818738], [171.42857142857142, 0.12771306118557813], [57.142857142857146, 0.0782649741917839], [64.51612903225806, 0.07135607527230461], [127.65957446808511, 0.065676126308472], [51.724137931034484, 0.04765587930734781], [68.96551724137932, 0.04586137274352434], [113.20754716981132, 0.032410353800951204], [250.0, 0.030369918652813157], [49.18032786885246, 0.029405569368220244], 'chords': [[0.0, 3.1, 'D#maj'], [3.1, 4.7, 'A#maj'], [4.7, 5.800000000000001, 'F:maj'], [5.800000000000001, 8.700000000000001, 'G:min'], [8.700000000000001, 11.0, 'F:maj'], [11.0, 13.200000000000001, 'D#maj'], [13.200000000000001, 14.3, 'G:min'], [14.3, 16.3, 'A#maj'], [16.3, 16.900000000000002, 'F:maj'], [16.900000000000002, 18.5, 'G:min'], [18.5, 19.1, 'A#maj'], [19.1, 19.900000000000002, 'D#maj'], [19.900000000000002, 20.400000000000002, 'A#maj'], [20.400000000000002, 21.900000000000002, 'C:maj'], [21.900000000000002, 22.8, 'G:min'], [22.8, 25.200000000000003, 'D#maj'], [25.200000000000003, 27.900000000000002, 'D#maj'], [27.900000000000002, 30.000000000000004, 'G:min']], 'downbeats': [[0.11, 1.0], [0.81, 2.0], [1.54, 3.0], [2.21, 4.0], [2.91, 1.0], [3.6, 2.0], [4.34, 3.0], [5.0, 4.0], [5.69, 1.0], [6.39, 2.0], [7.12, 3.0], [7.79, 4.0], [8.49, 1.0], [9.18, 2.0], [9.91, 3.0], [10.58, 4.0], [11.28, 1.0], [11.98, 2.0], [12.7, 3.0], [13.37, 4.0], [14.07, 1.0], [14.77, 2.0], [15.5, 3.0], [16.16, 4.0], [16.86, 1.0], [17.56, 2.0], [18.27, 3.0], [18.95, 4.0], [19.65, 1.0], [20.35, 2.0], [21.08, 3.0], [21.74, 4.0], [22.44, 1.0], [23.14, 2.0], [23.85, 3.0], [24.53, 4.0], [25.93, 2.0], [26.65, 3.0], [27.33, 4.0], [28.02, 1.0], [28.72, 2.0], [29.45, 3.0]], 'key': [[0.0026776420418173075, 0.5507075190544128, 0.000813949534669518, 0.0068492344580590725, 0.0008653984405100346, 0.04045381024479865, 0.015051022171974182, 0.002885610563805296, 0.00173998695107340681, 0.0084952414003579712, 0.0063829794496297836, 0.002278234576806426, 0.0019352855072915554, 0.0031427720095962286, 0.003523808353147507, 0.04947231337428093, 0.0013390464009717107, 0.014915909618139267, 0.003726801136508584, 0.002005897229537368, 0.0014139574486762285, 0.0013369121588766575, 0.2761196196079254, 0.001292902510613203]]}

Model: GPT-4o mini

Output: Unified music captions

The music clip unfolds with a brisk tempo, distinctively fast but balanced with a medium energy that fuels a palpable sense of urgency throughout its rhythm. This intensity carries a somber undertone, reflective in the music's low valence, leading listeners through a contemplative journey amid its forward momentum.

The harmonic structure pivots through a series of rich chords—beginning with the striking resonance of D#maj, transitioning through A#maj and F:maj, before settling into the darker reflections of G:min. This progression adds depth, embodying a landscape that evokes both movement and introspection. The interchange of these chords facilitates a pivotal emotional exploration, enhancing the track's melancholic essence while maintaining a danceable character, where reflection subtly intersperses with kinetic energy.

Synchronization of downbeats adds to the sense of drive, propelling the listener forward while imbuing each moment with a tender fragility. It creates a nuanced interplay between urgency and contemplation, inviting engagement without overwhelming the senses. In its entirety, the composition weaves a narrative that balances the frenetic pace of existence with the weight of emotional contemplation, crafting an immersive atmosphere where movement and melancholy coexist harmoniously.

Figure 3: An example of music caption, low-level features and generated unified music captions.

Input: (1) Video captions, (2) Unified Music Captions, (3) Lyrics understanding

Model: GPT-4o mini

Output: Video description

Overview

A mesmerizing journey through duality and introspection unfolds in this surreal music video. The brisk tempo mirrors the energetic yet contemplative soundscape, navigating between rich harmonies and softer reflective moments. The thoughtful interplay of organic and geometric visuals, complemented by shifting color tones, creates an immersive experience that invites deeper engagement with themes of transformation and optimism amidst chaos.

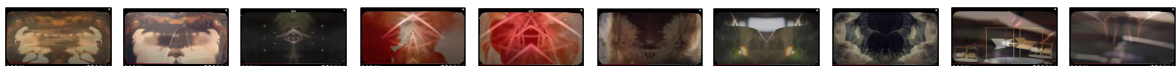
Frame-by-frame Breakdown

1. ****Blurred Reflection****: Begin with a soft, abstract blur, setting an enigmatic tone that hints at duality.
2. ****Obscured Humanity****: Introduce a warm, human form partially hidden, symbolizing vulnerability and connection.
3. ****Geometric Chaos****: Shift to darker visuals with floating geometric shapes, suggesting tension and direction.
4. ****Radiant Lines****: Bright red lines burst forth, enhancing the dynamic energy of the visual progression.
5. ****Intricate Motion****: More complex red layers emerge, indicating escalating movement and urgency.
6. ****Dreamlike Nature****: Transition to mirrored greenery, evoking serenity and a connection to the Earth.
7. ****Skyward Contemplation****: Move towards the calmness of a mirrored sky, fostering reflection and peace.
8. ****Whimsical Objects****: Introduce artistic mobile-like structures in soft tones, providing a sense of wonder.
9. ****Intimate Textures****: Zoom in on reflective surfaces, deepening the intimate atmosphere within the narrative.
10. ****Layered Complexity****: Explore deeper reflections and shadows, suggesting hidden depths beneath the surface.
11. ****Earthy Symmetry****: Return to natural settings with earthy tones, enhancing organic beauty and exploration.
12. ****Grounding Focus****: Highlight rich textures of the ground, reinforcing connection to nature.
13. ****Endless Horizons****: Present a symmetrical sky view, connoting endless possibilities through mirroring.
14. ****Tranquil Warmth****: A warm horizon emphasizes a feeling of openness and tranquility.
15. ****Reflective Closure****: Conclude with a blend of organic shapes and warm tones, inviting contemplation on the journey within the video.

This concise frame-by-frame breakdown supports a nuanced visualization that evokes the music's balance of urgency and reflection, mirroring the overarching themes expressed through lyrics.

Figure 4: An example of a full MV description.

• Original frames from MV:



• Frames generated by Text2video-Zero based on the MV description:



Figure 5: Frames from the original MV and generated by the Text2Video-Zero (Khachatryan et al., 2023) model.