
— Appendix —

MiraData: A Large-Scale Video Dataset with Long Durations and Structured Captions

Xuan Ju^{1,2*}, Yiming Gao^{1*}, Zhaoyang Zhang^{1†*}, Ziyang Yuan¹, Xintao Wang¹,
Ailing Zeng², Yu Xiong², Qiang Xu², Ying Shan¹
<https://github.com/mira-space/MiraData>

1 In the appendix, we first give out more details about the collection, selection, and annotation of
2 *MiraData* in Sec. A. Then, we provide additional experiment results of quantitative comparison,
3 qualitative comparison, and ablation in Sec. B. In Sec. C, we further explain the limitations, societal
4 impact, ethical issues and broad impact of our dataset. Finally, in Sec. D, we provide data acquisition,
5 data documentation, and data license for ease of data use.

6 A MiraData: Additional Details

7 A.1 Data Collection

8 We provide additional details about collecting YouTube video channels in this section. We select 7
9 categories that contain more rich motion and long video clips: (1) 3D engine-rendered scenes, (2)
10 city/scenic tours, (3) movies, (4) first-person perspective camera videos, (5) object creation/physical
11 law demonstrations, (6) timelapse videos, and (7) videos showcasing human motion.

12 The reason we choose these channels is as follows:

- 13 (1) For **3D engine-rendered scenes**, the videos are typically recorded in 3D rendering en-
14 gines with predefined physics laws. Thus, they often contain rich scene and perspective
15 changes, with relatively long continuous shots, making them suitable for learning long video
16 generation.
- 17 (2) **City/scenic tours** are usually filmed by people walking with handheld cameras in urban
18 or scenic areas. Consequently, the scenes are relatively continuous and possess strong 3D
19 spatial descriptive capabilities.
- 20 (3) **Movies** usually contain high-quality visuals and seamless transitions in the same scene,
21 allowing for a more comprehensive description of the same scene from different angles.
- 22 (4) **First-person-perspective camera videos** provide a perspective from the vantage point of
23 the person or device capturing the footage. Compared to city/scenic tours, this category
24 focuses more on extreme sports and typically uses camera lenses with slight distortion,
25 which offers a view from the eyes of the subject.
- 26 (5) **Object creation/physical law demonstration** often includes demonstrative videos focused
27 on a single perspective, such as baking tutorials or explanations of physical principles. Due
28 to their relatively simple scenes and clear procedural steps, these videos are beneficial for
29 learning physical laws in long videos.
- 30 (6) **The timelapse videos** capture a sequence of images at set intervals to record changes that
31 take place slowly over time, which represent processes that would be too slow to observe in

*Equal contribution. † Project Lead. ¹ARC Lab, Tencent PCG. ²The Chinese University of Hong Kong.

32 real-time. This would be helpful for the video generation model to learn real-world physics
 33 knowledge as indicated by MagicTime [1].
 34 (7) **Human motion videos** show human movements, such as speeches, dances, and model stage
 35 performances. Including this category will be beneficial for generating long videos that
 36 include localized limb movements of the human body. In Fig. 1, we provided two examples
 37 from each category to illustrate the differences between the various types of videos.

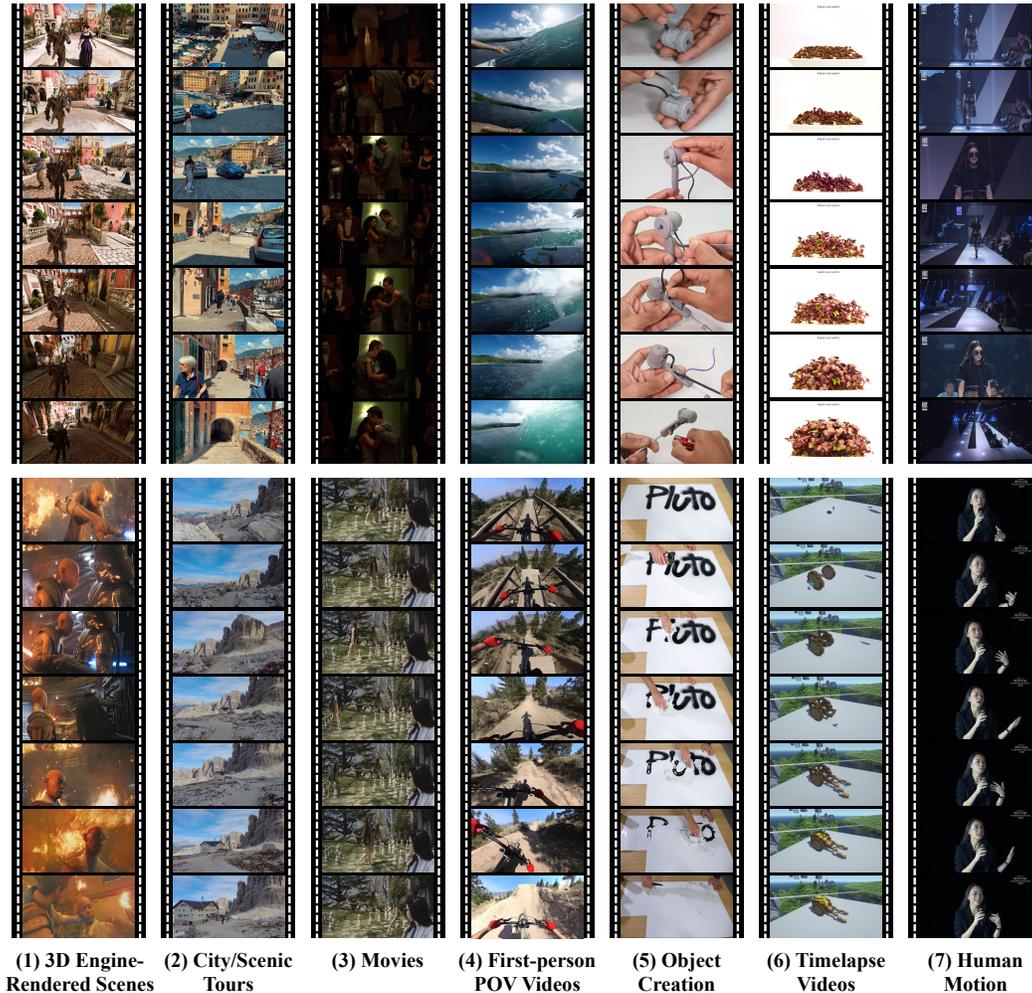


Figure 1: Video Examples From Each Category.

38 A.2 Video Splitting and Stitching

39 For video splitting, we use PySceneDetect² content-aware detection with a threshold of 26. This
 40 process may result in some incorrect separations when cutting long videos into small clips. To address
 41 this issue, we consider both content-coherent video transitions and wrong cuts.

42 To connect content-coherent video clips, we employ Qwen-VL-Chat[2] and LLaVA[3]. For each
 43 pair of adjacent video clips, we extract the 5th frame from the end of the former video and the 5th
 44 frame from the beginning of the latter video. These two frames are concatenated and input into the
 45 language models with the following prompt:

²<https://www.scenedetect.com/>

46 “Given two images shown on the right and the left, please determine whether the two images are
 47 similar to each other and coming from the same video. Please answer ‘Yes’ or ‘No’. You can start by
 48 examining the visual content of the two images. Look for similarities in various aspects, such as main
 49 objects, backgrounds, colors, lighting conditions, and spatial arrangements. Consider both global and
 50 local features within the images. For example, you should output ‘Yes’ for two images from different
 51 views of a scene. You should output ‘No’ for two unrelated images.”

52 The language models will output “Yes” or “No” as the answer. The two adjacent clips will be
 53 connected only when both models output “Yes”. To connect wrong cuts, we use ImageBind[4] and
 54 DINOv2[5] with thresholds of 0.6 and 0.85, respectively.

55 A.3 Video Selection

56 We list the filtering criteria in Tab. 1. Average Optical flow measures the overall motion across the
 57 video sequence, giving an idea of how much movement is occurring on average. Image Max 30%
 58 Optical Flow identifies each image’s maximum 30% optical flow values. This can focus on the part
 59 of the image that contains the largest movement. Temporal Min 30% Optical Flow the minimum
 60 optical flow values within the bottom 30% of frames in terms of motion intensity, giving the results
 61 of the least dynamic parts of the image sequence by focusing on the frames with the lowest motion.
 62 The Average Aesthetic Score is assessed using the Laion-Aesthetic[6] Aesthetic Score Predictor and
 63 averaged among frames. Average Color is the average of the color of every pixel in frames. Temporal
 64 Max 80% Color identifies the maximum color values within the top 80% of frames, which is the
 65 brightest. Temporal Min 80% Color identifies the maximum color values within the bottom 80% of
 66 frames, which is the darkest. Contain NSFW identifies whether the frames contain NSFW content.

Table 1: **Filtering Criteria of MiraData.** We offer five versions of MiraData, each filtered using different criteria to cater to various research needs and preferences.

Metrics	788K Version	330K Version	93K Version	42K Version	9K Version
Average Optical Flow	-	>2.0	>3.0	>4.0	>4.5
Image Max 30% Optical Flow	-	-	>4.3	>4.8	>5.1
Temporal Min 30% Optical Flow	-	-	>2.5	>3.5	>4.0
Average Aesthetic Score	-	-	>3.0	>5.0	>5.5
Average Color	-	>25.0 <230.0	>25.0 <230.0	>35.0 <220.0	>35.0 <220.0
Temporal Max 80% Color	-	-	<235.0	<225.0	<225.0
Temporal Min 80% Color	-	-	>20.0	>30.0	>30.0
Contain NSFW	No	No	No	No	No

67 A.4 Video Caption

68 To facilitate the comprehension of videos by GPT-4V, we extract eight uniformly sampled frames
 69 from each video, arranging them in a 2×4 grid within a single image. Alongside this 2×4 grid image,
 70 we meticulously design a prompt to enable GPT-4V to perceive this image as a video thumbnail.
 71 Following DALL-E3 [7], we bias GPT-4V to yield video descriptions conducive to the learning of a
 72 text-to-video generation model. Our initial step utilizes Panda-70M [8] to produce a “short caption”
 73 that delineates the primary subject and actions, serving as an additional cue for GPT-4V. Specifically,
 74 our prompt begins with the following guiding content:

75 A wide image is given containing a 2×4 grid of 8 equally spaced video frames. They’re arranged
 chronologically from left to right, and then from top to down, all separated by white borders. This
 video depicts “*Short Captions*”. Please imagine the video based on the sequence of 8 frames, and
 provide a faithfully concise description of the following content:

76 We further instruct GPT-4V to generate dense descriptions of videos. In addition, we introduce
 77 structured captions to obtain more intricate information. To procure more precise, detailed, and
 78 fine-grained structured captions, we carefully craft prompts that inquire about various aspects of
 79 the video, including the main object, background, camera movement, and video style. The specific
 80 prompts are described below:

98 improvements in motion strength while maintaining temporal and 3D consistency compared to the
 99 model trained on WebVid-10M.

Table 2: **Quantitative Comparison of MiraDiT trained on MiraData and WebVid-10M [9].** \uparrow and \downarrow means higher/lower is better. 1) - 14) are metrics of MiraBench, where DD for Dynamic Degree, TS for Tracking Strength, DTC for DINO Temporal Consistency, CTC for CLIP Temporal Consistency, TMS for Temporal Motion Smoothness, MAE for Mean Absolute Error, RMSE for Root Mean Square Error, AQ for Aesthetic Quality, IQ for Imaging Quality, CA for Camera Alignment, MOA for Main Object Alignment, BA for Background Alignment, SA for Style Alignment, and OA for Overall Alignment.

Metrics	Temporal Motion Strength		Temporal Consistency			3D Consistency	
	1) DD \uparrow	2) TS \uparrow	3) DTC \uparrow	4) CTC \uparrow	5) TMS \uparrow	6) MAE $\downarrow \times 10^{-2}$	7) RMSE $\downarrow \times 10^{-1}$
OpenSora [11]	7.65	16.07	12.34	13.20	13.70	75.45	10.39
	7.48	15.21	11.86	12.94	13.03	78.23	11.06
	7.59	15.78	12.01	13.04	13.42	77.18	10.82
VideoCrafter2 [10]	1.71	6.72	6.41	6.36	6.60	101.55	13.05
	3.01	8.52	9.12	8.89	9.23	120.05	15.13
	2.02	6.91	6.53	6.42	6.84	99.84	12.87
MiraDiT (WebVid-10M [9])	7.12	22.36	20.24	20.97	21.86	91.48	12.11
	6.93	21.74	20.23	20.49	22.30	90.11	11.96
	7.18	22.52	20.30	20.99	21.95	91.31	12.12
MiraDiT (<i>MiraData</i>)	15.46	49.47	43.78	45.95	47.24	85.27	11.74
	15.32	49.41	43.66	45.85	47.19	84.22	11.68
	16.03	50.26	44.01	45.99	47.32	86.11	11.94

Metrics	Visual Quality		Text-Video Alignment				
	8) AQ $\uparrow \times 10^{-1}$	9) IQ \uparrow	10) CA \uparrow	11) MOA \uparrow	12) BA \uparrow	13) SA \uparrow	14) OA \uparrow
OpenSora [11]	47.10	59.54	12.40	18.12	13.20	13.35	16.12
	44.28	60.14	12.38	17.93	13.41	13.39	16.82
	48.01	58.39	12.01	18.25	13.38	13.96	16.57
VideoCrafter2 [10]	58.69	64.96	12.00	17.90	11.25	12.15	16.90
	57.96	64.86	12.09	17.86	11.63	12.09	16.59
	58.28	64.99	11.97	17.78	11.42	12.04	16.68
MiraDiT (WebVid-10M [9])	43.11	58.58	12.35	14.32	11.90	12.32	15.31
	43.01	57.74	12.43	14.29	12.01	12.29	16.01
	43.42	59.00	12.42	14.33	11.96	12.21	15.48
MiraDiT (<i>MiraData</i>)	49.90	63.71	12.66	14.67	12.18	12.59	16.66
	50.21	63.58	12.86	14.69	12.25	12.53	16.84
	49.53	63.70	12.73	14.69	12.17	12.64	16.73

100 To assess MiraDiT’s performance on other benchmark datasets, we test the performance of MiraDiT
 101 on the recent text-to-video benchmark, T2V-CompBench [12]. T2V-CompBench includes 7 metrics
 102 designed to evaluate the alignment of generated videos with the corresponding text prompts: (1)
 103 Consistent Attribute Binding: Evaluates whether object attributes remain consistent throughout the
 104 generated video frames. (2) Dynamic Attribute Binding: Assesses if the generated video accurately
 105 reflects changes in object attributes. (3) Spatial Relationship: Determines if the generated video
 106 adheres to the spatial relationships specified in the text prompt. (4) Motion Binding: Assesses the
 107 correctness of the object’s motion direction in the generated video. (5) Action Binding: Evaluates
 108 the accuracy of the object action categories in the generated video. (6) Object Interactions: Tests
 109 the model’s ability to generate dynamic interactions between objects. (7) Generative Numeracy:
 110 Evaluates the accuracy in the number of objects generated as specified in the text prompt. Results
 111 show that MiraDiT trained on MiraData achieves much better results on all metrics compare to that
 112 trained on WebVid-10M. Moreover, MiraDiT trained on MiraData have the best results on Dynamic
 113 Attribute Binding, further illustrates the advantages of training with high-dynamic, detailed-captioned
 114 data. MiraDiT trained on MiraData also achieves a relatively advanced results in all open-source
 115 text-to-video generation models. However, we must point out, that this comparison is unfair, as
 116 different models were trained using different computational resources and distinct models, making it
 117 impossible to assess the quality of MiraData relative to other training data. Moreover, the evaluation
 118 prompts in T2V-CompBench primarily consist of short captions with only a single simple sentence,
 119 which limits MiraData’s ability to fully showcase its strengths.

Table 3: **T2V-CompBench** evaluation results of MiraDiT trained on MiraData and WebVid-10M. Best results are shown in bold.

Method	Consist-attnr \uparrow	Dynamic-attnr \uparrow	Spatial \uparrow	Motion \uparrow	Action \uparrow	Interaction \uparrow	Numeracy \uparrow
ModelScope	0.5483	0.1654	0.4220	0.2552	0.4880	0.7075	0.2066
ZeroScope	0.4495	0.1086	0.4073	0.2319	0.4620	0.5550	0.2378
Latte	0.5325	0.1598	0.4476	0.2187	0.5200	0.6625	0.2187
Show-1	0.6388	0.1828	0.4649	0.2316	0.4940	0.7700	0.1644
VideoCrafter2	0.6750	0.1850	0.4891	0.2233	0.5800	0.7600	0.2041
Open-Sora 1.1	0.6370	0.1762	0.5671	0.2317	0.5480	0.7625	0.2363
Open-Sora 1.2	0.6600	0.1714	0.5406	0.2388	0.5717	0.7400	0.2556
Open-Sora-Plan v1.0.0	0.5088	0.1562	0.4481	0.2147	0.5120	0.6275	0.1650
Open-Sora-Plan v1.1.0	0.7413	0.1770	0.5587	0.2187	0.6780	0.7275	0.2928
AnimateDiff	0.4883	0.1764	0.3883	0.2236	0.4140	0.6550	0.0884
VideoTetris	0.7125	0.2066	0.5148	0.2204	0.5280	0.7600	0.2609
LVD	0.5595	0.1499	0.5469	0.2699	0.4960	0.6100	0.0991
MagicTime	-	0.1834	-	-	-	-	-
MiraDiT (WebVid-10M)	0.6012	0.1972	0.4438	0.2250	0.5156	0.6075	0.1909
MiraDiT (MiraData)	0.6825	0.2302	0.4622	0.2321	0.6340	0.7373	0.2234

120 B.3 Role of Video Duration

121 To evaluate the effectiveness of *MiraData* on long-duration video generation, we train a dynamic
 122 frame rate video generation model that supports arbitrary length video generation from 0 to 20s on
 123 *MiraData* and WebVid respectively. Tab. 4 presents the results for 5s, 10s, and 20s videos. Tab. 4
 124 presents the results for 5s, 10s, and 20s videos. The experimental results demonstrate that our
 125 *MiraData* achieves significantly better motion strength and dynamic degree compared to the model
 126 trained on WebVid-10M, while maintaining consistent temporal and 3D consistency. Furthermore,
 127 *MiraData* yields higher aesthetic scores, attributed to its high video visual quality (e.g., resolution
 128 and aesthetic scores). As the generated video duration increases, *MiraData*’s performance in motion
 129 intensity and aesthetic scores improves, benefiting from the longer video clips in our dataset.

Table 4: **Ablation on Video Duration.** \uparrow and \downarrow means higher/lower is better. 1) - 14) are metrics of MiraBench. Refer to Tab. 2 for a detailed explanation of annotation.

Metrics	Temporal Motion Strength		Temporal Consistency			3D Consistency		
	1) DD \uparrow	2) TS \uparrow	3) DTC \uparrow	4) CTC \uparrow	5) TMS \uparrow	6) MAE $\downarrow \times 10^{-2}$	7) RMSE $\downarrow \times 10^{-1}$	
WebVid-10M [9]	5s	7.12	22.36	20.24	20.97	21.86	91.48	12.11
	10s	4.82	24.99	23.23	23.63	24.62	94.62	12.53
	20s	4.73	63.74	57.18	59.06	62.33	99.62	13.01
<i>MiraData</i>	5s	15.46	49.47	43.78	45.95	47.24	85.27	11.74
	10s	5.23	27.06	25.22	25.67	26.55	89.44	12.08
	20s	6.41	84.41	76.19	78.61	82.48	96.66	12.94

Metrics	Visual Quality		Text-Video Alignment					
	8) AQ $\uparrow \times 10^{-1}$	9) IQ \uparrow	10) CA \uparrow	11) MOA \uparrow	12) BA \uparrow	13) SA \uparrow	14) OA \uparrow	
WebVid-10M [9]	5s	43.11	58.58	12.35	14.32	11.90	12.32	15.31
	10s	40.98	59.60	0.12	12.99	11.61	11.91	13.65
	20s	37.93	59.11	12.07	12.32	11.92	11.48	13.31
<i>MiraData</i>	5s	49.90	63.71	12.66	14.67	12.18	12.59	16.66
	10s	42.60	61.47	11.97	13.62	11.17	11.77	14.94
	20s	40.36	59.32	12.00	13.96	11.09	11.69	14.56

130 C Limitations and Potential Negative Societal Impacts

131 C.1 Limitations and Future Work

132 Despite the advancements and contributions of our work, several limitations need to be acknowledged
 133 and addressed in future research:

- 134 • **Dataset Diversity and Coverage.** Although *MiraData* presents a substantial improvement
 135 over existing datasets, it may still lack comprehensive diversity in terms of content, genre,

136 and cultural representation. The dataset’s reliance on manually selected channels might
137 introduce a bias towards certain types of videos, potentially affecting the generalizability of
138 models trained on it. Future work could focus on expanding the dataset by including a wider
139 variety of sources and a more balanced representation of different content types.

140 • **Scalability of the Data Curation Pipeline.** The current data curation pipeline, while
141 effective, might face challenges in scalability, particularly in handling the growing volume
142 of video content and the complexity of annotations required for maintaining high quality.
143 Automating more aspects of the data curation process with more efficient machine learning
144 models could improve scalability.

145 • **Caption Quality.** While the structured captions in *MiraData* are more detailed than previous
146 datasets, there might still be instances where the captions do not fully capture the nuanced
147 details of the video content. Additionally, using automated captioning tools, despite their
148 high accuracy, can occasionally result in errors or ambiguities. Enhancing the captioning
149 process by integrating human-in-the-loop methods can improve the quality and accuracy
150 of captions. Furthermore, iterative refinement of captions based on feedback from domain
151 experts and end-users could help in generating more precise and informative descriptions.

152 • **Evaluation Metrics.** The proposed *MiraBench* benchmark, although comprehensive, may
153 not fully cover all aspects of video generation quality, especially those related to subjective
154 human perceptions such as creative quality. Incorporating human evaluations alongside
155 automated metrics can provide a more holistic assessment of generated videos.

156 C.2 Potential Negative Societal Impacts and Solutions

157 The construction of video datasets can lead to possible negative societal impacts such as: (1)
158 Misinformation and Deepfakes. Advances in text-to-video generation models like Sora, particularly
159 those that produce highly realistic and detailed videos, raise significant concerns about the potential
160 for creating deepfakes. These realistic fake videos can be used to spread misinformation, manipulate
161 public opinion, or damage reputations. To solve this, we need to implement robust detection
162 mechanisms and watermarking-generated content to help identify and prevent the misuse of AI-
163 generated videos. Additionally, establishing ethical guidelines and legal frameworks to regulate the
164 use of such technology is crucial. (2) Including Personally Identifiable Information. Collecting videos
165 from various platforms could result in the inclusion of content that contains personally identifiable
166 information, such as faces, locations, or other identifiable features, without consent. We should
167 implement stringent data anonymization techniques and manual review processes to ensure that any
168 PII is either removed or consent is obtained before including such data in the dataset. (3) Bias and
169 Stereotyping. If the dataset used for training contains biased or stereotypical representations, the
170 generated content may perpetuate these biases, leading to harmful societal stereotypes and reinforcing
171 negative perceptions. So, we need to actively curate a diverse and balanced dataset that represents
172 various demographics and perspectives, which can help mitigate bias. Regularly auditing the models
173 for biased outputs and retraining them on more balanced datasets can further reduce this risk.

174 D Data Acquisition and License

175 **Data Acquisition.** Data downloading link is: <https://github.com/mira-space/MiraData>.

176 **License.** This dataset is made available for informational purposes only. No license, whether implied
177 or otherwise, is granted in or to such dataset (including any rights to copy, modify, publish, distribute
178 and/or commercialize such dataset), unless you have entered into a separate agreement for such
179 rights. Such dataset is provided as-is, without warranty of any kind, express or implied, including any
180 warranties of merchantability, title, fitness for a particular purpose, non-infringement, or that such
181 dataset is free of defects, errors or viruses. In no event will our institution be liable for any damages
182 or losses of any kind arising from the dataset or your use thereof.

References

- 183
- 184 [1] S. Yuan, J. Huang, Y. Shi, Y. Xu, R. Zhu, B. Lin, X. Cheng, L. Yuan, and J. Luo, “Mag-
185 ictime: Time-lapse video generation models as metamorphic simulators,” *arXiv preprint*
186 *arXiv:2404.05014*, 2024.
- 187 [2] J. Bai, S. Bai, S. Yang, S. Wang, S. Tan, P. Wang, J. Lin, C. Zhou, and J. Zhou, “Qwen-vl: A
188 frontier large vision-language model with versatile abilities,” *arXiv preprint arXiv:2308.12966*,
189 2023.
- 190 [3] H. Liu, C. Li, Y. Li, and Y. J. Lee, “Improved baselines with visual instruction tuning,” 2023.
- 191 [4] R. Girdhar, A. El-Nouby, Z. Liu, M. Singh, K. V. Alwala, A. Joulin, and I. Misra, “Imagebind:
192 One embedding space to bind them all,” in *Proceedings of the IEEE/CVF Conference on*
193 *Computer Vision and Pattern Recognition*, pp. 15180–15190, 2023.
- 194 [5] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza,
195 F. Massa, A. El-Nouby, *et al.*, “Dinov2: Learning robust visual features without supervision,”
196 *arXiv preprint arXiv:2304.07193*, 2023.
- 197 [6] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes,
198 A. Katta, C. Mullis, M. Wortsman, *et al.*, “Laion-5b: An open large-scale dataset for training
199 next generation image-text models,” *Advances in Neural Information Processing Systems*,
200 vol. 35, pp. 25278–25294, 2022.
- 201 [7] J. Betker, G. Goh, L. Jing, T. Brooks, J. Wang, L. Li, L. Ouyang, J. Zhuang, J. Lee, Y. Guo,
202 *et al.*, “Improving image generation with better captions,” *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, vol. 2, no. 3, p. 8, 2023.
- 204 [8] T.-S. Chen, A. Siarohin, W. Menapace, E. Deyneka, H.-w. Chao, B. E. Jeon, Y. Fang, H.-Y. Lee,
205 J. Ren, M.-H. Yang, *et al.*, “Panda-70m: Captioning 70m videos with multiple cross-modality
206 teachers,” *arXiv preprint arXiv:2402.19479*, 2024.
- 207 [9] M. Bain, A. Nagrani, G. Varol, and A. Zisserman, “Frozen in time: A joint video and image
208 encoder for end-to-end retrieval,” in *Proceedings of the IEEE/CVF International Conference on*
209 *Computer Vision*, pp. 1728–1738, 2021.
- 210 [10] H. Chen, Y. Zhang, X. Cun, M. Xia, X. Wang, C. Weng, and Y. Shan, “Videocrafter2: Overcom-
211 ing data limitations for high-quality video diffusion models,” 2024.
- 212 [11] Z. Zangwei, P. Xiangyu, L. Shenggui, L. Hongxing, Z. Yukun, L. Tianyi, P. Xiangyu, Z. Zangwei,
213 S. Chenhui, Y. Tom, W. Junjie, and Y. Chenfeng, “Opensora,” 2024.
- 214 [12] K. Sun, K. Huang, X. Liu, Y. Wu, Z. Xu, Z. Li, and X. Liu, “T2v-compbench: A comprehensive
215 benchmark for compositional text-to-video generation,” *arXiv preprint arXiv:2407.14505*, 2024.