

MANTISSCORE: A Reliable Fine-grained Metric for Video Generation

Anonymous ACL submission

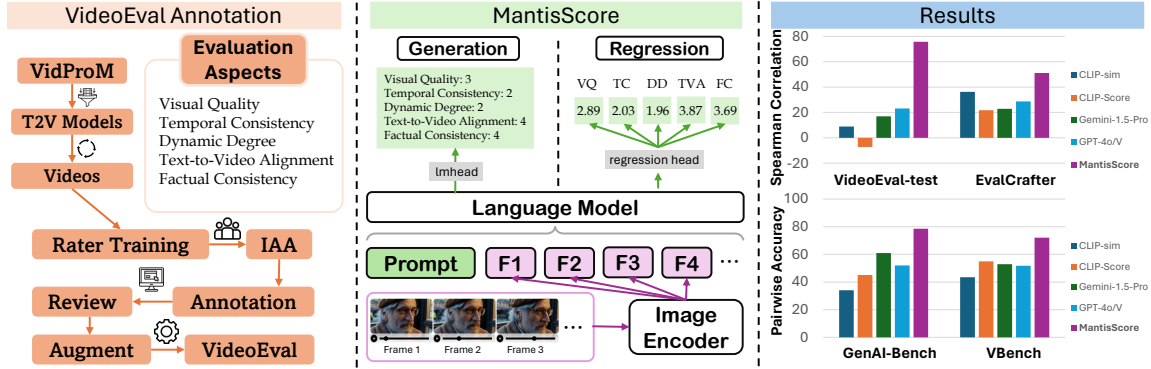


Figure 1: Construction process of VIDEOEVAL dataset and illustration of MANTISSCORE.

Abstract

The recent years have witnessed great advances in text-to-video generation. However, the video evaluation metrics have lagged significantly behind, which fails to produce an accurate and holistic measure of the generated videos’ quality. The main barrier is the lack of high-quality human rating data. In this paper, we release VIDEOEVAL, the first large-scale multi-aspect video evaluation dataset. VIDEOEVAL consists of high-quality human-provided ratings for 5 video evaluation aspects on the 37.6K videos generated from 11 existing popular video generative models. We train MANTISSCORE based on VIDEOEVAL to enable automatic video quality assessment. Experiments show that the Spearman correlation between MANTISSCORE and humans can reach 77.1 on VIDEOEVAL-test, beating the prior best metrics by about 50 points. Further result on the held-out EvalCrafter, GenAI-Bench, and VBench, show that MANTISSCORE is highly generalizable and still beating the prior best metrics by a remarkable margin. We observe that using Mantis as the based model consistently beats that using Idefics2 and VideoLLaVA, and the regression-based model can achieve better results than the generative ones. Due to its high reliability, we believe MANTISSCORE can serve as a valuable tool for accelerate video generation research.

1 Introduction

Powerful text-to-video (T2V) generative models have been exponentially emerging these days. In

2023 and 2024, we have witnessed an array of T2V models like Sora (OpenAI, 2024b), Runway Gen-2 (Esser et al., 2023), Lumiere (Bar-Tal et al., 2024), Pika¹, Luma-AI², Kling³, Emu-video (Girdhar et al., 2023), StableVideoDiffusion (Blattmann et al., 2023a). These models have shown their potential to generate longer-duration, higher-quality, and more natural videos. Despite significant advancements in video generation models, the evaluation metrics of video generation is lagging behind.

The recent literature has adopted a wide range of metrics to do video quality assessments. However, these metrics suffer from the following issues: (1) they can only be used to evaluate visual quality or aesthetics, while failing to capture aspects like motion smoothness, factual consistency, etc. Examples of such metrics include CLIP (Radford et al., 2021b), DINO (Caron et al., 2021), BRISQUE (Mittal et al., 2012a), FVD (Unterthiner et al., 2019), and IS (Salimans et al., 2016). (2) some metrics focus only on a single mean opinion score (MOS), failing to provide fine-grained subscores across different multiple aspects. Examples include T2VQA (Kou et al., 2024b), FastVQA (Wu et al., 2022), and DOVER (Wu et al., 2023). Several works (Ku et al., 2023; Bansal et al., 2024) propose to prompt multi-modal large-language-models (MLLM) like GPT-4o (Achiam et al., 2023) or

¹<https://pika.art/home>

²<https://lumalabs.ai/dream-machine>

³<https://kling.kuaishou.com/>

Gemini-1.5 (Reid et al., 2024) to produce multi-aspect quality assessment for given videos. However, our experiments show that they also have low correlation with humans.

The biggest barrier to build reliable video metrics is the lack of high-quality human-annotated dataset. To overcome this barrier, we curate VIDEOEVAL, the first large-scale, multi-aspect video evaluation dataset. We select prompts from VidProM (Wang and Yang, 2024), and use 11 popular text-to-video models, including Pika, Lavie (Wang et al., 2023c), SVD (Blattmann et al., 2023a), etc, to generate videos of various quality based on these prompts. We define five key aspects for evaluation in Table 2, and each aspect is scored from 1 (bad) to 4 (perfect). For annotation, we trained 20 raters to perform a multi-aspect rating over individual generated videos. We have collected ratings for a total of 37.6K videos. We iterate multiple rounds of refinement to ensure a high inter-annotation-agreement (IAA) ratio over 60% for all five aspects.

To build the video evaluator, we select Mantis-Idefics2-8B (Jiang et al., 2024a) as our main backbone model due to its superior ability to handle multi-image and video content, accommodating up to 128 video frames and supporting native resolution. After fine-tuning Mantis on VIDEOEVAL-train, we get our video evaluator, MANTISSCORE. Experiments show that we achieve a Spearman correlation of 77.1 on VIDEOEVAL-test and 59.5 on EvalCrafter (Liu et al., 2023b) for the text-to-video alignment aspect, surpassing the best baseline by 54.1 and 4.4 respectively. The pairwise comparison accuracy gets 78.5 on GenAI-Bench (Jiang et al., 2024b) video preference part, and 72.1 in average on 5 aspects of VBench (Huang et al., 2023), surpassing the previous best baseline by 11.4 and 9.6 respectively. Additional ablation studies with different backbone models confirmed that the Mantis-based metric provides a gain of 12.1 compared to using the Idefics2-based metric. Due to the significant improvement, we believe that MANTISSCORE can serve as the reliable metrics for future video generative models.

2 Related Work

2.1 Text-to-Video Generative Models

Recent progress in diffusion models (Ho et al., 2020; Rombach et al., 2022) has significantly pushed forward the development of Text-to-Video

(T2V) generation. Given a text prompt, the T2V generative model can synthesize new video sequences that didn't previously exist (Wang et al., 2023c; OpenAI, 2024b; Chen et al., 2023a, 2024a; Henschel et al., 2024; Bar-Tal et al., 2024). Early diffusion-based video models generally build upon Text-to-Image (T2I) models, adding a temporal module to extend itself into the video domain (Wang et al., 2023c; Chen et al., 2023c). Recent T2V generation models are directly trained on videos from scratch. Among these, models based on Latent Diffusion Models (LDMs) have gained particular attention for their effectiveness and efficiency (Zhou et al., 2022; An et al., 2023; Blattmann et al., 2023b). While the other works used the pixel-based Diffusion Transformers (DiT) also achieve quality results (Gupta et al., 2023; Chen et al., 2023b; OpenAI, 2024b).

2.2 Video Quality Assessment

As the current progress of Text-to-Video generative models leaves it uncertain how close we are to reaching the objective, researchers have worked on evaluation methods to benchmark the generative models. Common methods involve the use of FVD (Unterthiner et al., 2018) and CLIP (Radford et al., 2021a) to evaluate the quality of frames and the text-frames alignment respectively. However, other aspects like subject consistency, temporal consistency, factualness cannot be captured by these metrics. Recent works like VBench (Huang et al., 2023) proposes to use different DINO (Caron et al., 2021), optical flow (Horn and Schunck, 1981) to reflect these aspects. However, the correlation with human judgment is relatively low. For example, most models have subject/background consistency scores over 97% in VBench, which is a massive overestimation of the current T2V models' true capability. Another work EvalCrafter (Liu et al., 2023b) instead resorts to human raters to perform comprehensive evaluation.

A recent work VideoPhy (Bansal et al., 2024) follows VIEScore (Ku et al., 2023) prompt large multi-modal models like Gemini (Reid et al., 2024) and GPT-4o (Achiam et al., 2023) to provide quality assessment. However, our later study shows that these multimodal language models also achieve very low agreement with human raters. A concurrent work T2VQA (Kou et al., 2024a) also proposes to train a quality assessment model on human-annotated video ratings. However, there are a few distinctions. Firstly, our dataset contains ratings

for multiple aspects. Secondly, our dataset is 4x larger than the T2VQA dataset. Thirdly, our metric is built on pre-trained video-language foundation models to maximize its performance.

3 VIDEOEVAL

This section introduces the construction process of our dataset, VIDEOEVAL, for training video evaluators. We start by explaining how we gathered and filtered diverse text prompts for video generation, followed by the video-generation processes using 11 selected text-to-video models. Next, we outline the annotation pipeline that guides raters to score videos across multiple aspects defined in Table 2. We also include supplementary data to enhance robustness. Finally, we summarize the dataset statistics in Table 1, with 760 examples designated as the test set for evaluation.

3.1 Data preparation

Prompt Sources We utilize VidProM (Wang and Yang, 2024), a dataset containing extensive text-to-video pairs from different models. VidProM’s video-generation prompts are diverse and semantically rich, derived from real-world user inputs. To create a manageable subset from the 1.04 million unique prompts, we apply two filters: a length filter and an NSFW filter. The length filter eliminates prompts with fewer than 5 words or more than 100 words. The NSFW filter removes prompts with a high probability of containing inappropriate content. After filtering, we perform random down-sampling to obtain a set of 44.5K prompts, 31.6K of them are used in video generation and some videos may have the same text prompt.

Video Generation We select 11 text-to-video (T2V) generative models (shown in Table 1) with various capabilities so that the quality of the generated video ranges from high to low in a balanced way. Some videos are pre-generated in the VidProM dataset, including Pika, Text2Video-Zero (Khachatryan et al., 2023), VideoCrafter2 (Chen et al., 2024a), and ModelScope (Wang et al., 2023a), whereas the others are generated by ourselves or collected from the Internet (i.e. SoRA). To eliminate differences between models in subsequent annotation stage, we normalize the videos into a unified format. First, we standardized the frame rate to 8 fps to address discrepancies in temporal consistency between high and low fps videos. Specifically, for high

frame rate model Pika and AnimateDiffusion (Guo et al., 2023) we use frame down sampling, while for low frame rate model like Text2Video-Zero, we employed frame interpolation (Huang et al., 2022) on it. Details are shown in Appendix E. Additionally, we cropped Pika videos to remove the watermark, making them indistinguishable from other models. Ultimately, we obtained 33.6K videos from 11 T2V models, along with their generation prompts.

3.2 Annotation Pipeline

Evaluation Dimensions As discussed in section 1, fine-grained and multi-aspect rating of videos is crucial for enhancing both the reliability and explainability of the video evaluator. Inspired by VBench (Huang et al., 2023) and EvalCrafter (Liu et al., 2023b), and FETV (Liu et al., 2023c), we propose five key dimensions for text-to-video evaluation, detailed in Table 2. These dimensions encompass both low-level vision aspects, such as Visual Quality, which evaluates basic visual impressions, and higher-level aspects, like Text-to-Video Alignment and Factual Consistency, which require a deep understanding of world knowledge, is a capability previous metrics do not have. Besides definition, a checklist for error points for each dimension is also provided to assist the rater in contributing more accurate and consistent rating. Detailed are provided in Table 8.

Annotation We hired 20 expert raters, with each rater performing rating for 1K-2K videos. Our raters are mostly college graduate students. For each aspect, there are three available ratings, 1 (Bad), 2 (Average), and 3 (Good), the score 4 (Perfect) is post-annotated, as described in the subsection 3.3. To ensure the consistency and quality of the annotations, we conducted a system training for each rater. Initially, we conducted a pilot training session with examples of multi-aspect ratings for various videos. Following this, multiple rounds of small-scale annotation were conducted to compute the inter-annotator agreement (IAA) across five aspects, as shown in Table 3. The results indicate a high score-matching ratio for all aspects, along with Fleiss’ κ (Fleiss and Cohen, 1973) and Krippendorff’s α (Krippendorff, 2011) metrics, with values around 0.4 or 0.5, suggesting sufficient agreement to proceed with large-scale annotation. The annotation process takes roughly 4 weeks to finish.

Base Model or Video Type	Video Source	Total Size	Resolution	Duration	FPS	Score
Human Annotated Videos						
Pika	VidProM	4.6k	(768, 480)	3.0s	8	[1-4]
Text2Video-Zero (Khachatryan et al., 2023)	VidProM	4.6k	(512, 512)	2.0s	8	[1-4]
VideoCrafter2 (Chen et al., 2024a)	VidProM	4.9k	(512, 320)	2.0s	8	[1-4]
ModelScope (Wang et al., 2023a)	VidProM	4.5k	(256, 256)	2.0s	8	[1-4]
LaVie-base (Wang et al., 2023c)	Generated	3.2k	(512, 320)	2.0s	8	[1-4]
AnimateDiff (Guo et al., 2023)	Generated	1.4k	(512, 512)	2.0s	8	[1-4]
LVDM (He et al., 2022)	Generated	3.1k	(256, 256)	2.0s	8	[1-4]
Hotshot-XL (Mullan et al., 2023)	Generated	3.2k	(512, 512)	1.0s	8	[1-4]
ZeroScope-576w (Sterling, 2024)	Generated	2.2k	(256, 256)	2.0s	8	[1-4]
Fast-SVD (Blattmann et al., 2023a)	Generated	1.0k	(1024, 576)	3.0s	8	[1-4]
SoRA-Clip (OpenAI, 2024b)	Collected	0.9k	various	2.0/3.0s	8	[1-4]
Augmented Videos						
DiDeMo (Hendricks et al., 2017)	Real	2.0k	various	2.0/3.0s	8	4
Panda70M (Chen et al., 2024b)	Real	2.0k	various	2.0/3.0s	8	4

Table 1: Statistics of our curated VIDEOEVAL for training video-generation evaluator. It consists of 33.6K human-scored videos across multiple aspects, with 4k real-world videos collected from DiDeMo (Hendricks et al., 2017) and Panda70M (Chen et al., 2024b) as the supplementary data. Ultimately, we get 37.6K high-quality rated videos as the final VIDEOEVAL.

Aspect	Definition
Visual Quality (VQ)	the quality of the video in terms of clearness, resolution, brightness, and color
Temporal Consistency (TC)	the consistency of objects or humans in video
Dynamic Degree (DD)	the degree of dynamic changes
Text-to-Video Alignment (TVA)	the alignment between the text prompt and the video content
Factual Consistency (FC)	the consistency of the video content with common-sense and factual knowledge

Table 2: The five evaluation aspects of VIDEOEVAL and their definitions.

IAA metric	VQ	TC	DD	TVA	FC
Trial 1 (#=30)					
Match Ratio	0.733	0.706	0.722	0.678	0.633
Kappa	0.369	0.414	0.413	0.490	0.265
Alpha	0.481	0.453	0.498	0.540	0.365
Trial 2 (#=100)					
Match Ratio	0.787	0.699	0.913	0.570	0.727
Kappa	0.088	0.562	0.565	0.125	-0.089
Alpha	0.078	0.579	0.620	0.205	-0.106

Table 3: Inter-Annotator Agreement (IAA) analysis results considering Matching Ratio, Fleiss’ κ , and Krippendorff’s α on the two trial annotations.

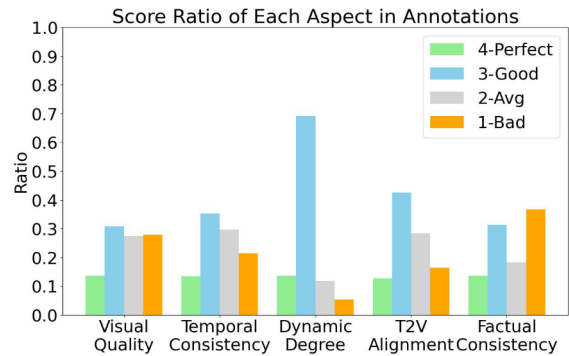


Figure 2: The rating distribution on all the videos.

Review We conduct random checks on human scores during the annotating process. Once we find the exceeded unqualified ratio in certain rater, we promptly communicate with the respective rater and review the annotations for that segment of the video. This helps calibrate the annotation provided by that rater during the relevant period. For example, we found several raters are too lenient and tend to give high scores to unqualified videos. We then step in to make sure they are aligned with our under-

standing of evaluation dimensions. With periodical random inspection on annotating, we completed the large-scale annotation of 33.6K videos and moved to the data augmentation stage.

3.3 Dataset Augmentation

To enhance the robustness of VIDEOEVAL dataset, we incorporated post-augmentation into the dataset. Firstly, expert raters will review the excellent videos (all aspects are scored 3) again to select perfect ones and raise their scoring to 4 (Perfect) in

280	certain aspects, particularly among the SoRA and	sine similarities of between adjacent frames	328
281	FastSVD (Blattmann et al., 2023a) videos.	features, following VBench (Huang et al.,	329
282	Additionally, we gather 4k real-world videos	2023). Additionally, we calculate SSIM be-	330
283	from the DiDeMo (Hendricks et al., 2017) and	tween adjacent frames, denoted as SSIM-sim.	331
284	Panda70M (Chen et al., 2024b) with each video		
285	accompanied by a text description. We select and	3. Dynamic Degree. We uniformly sample four	332
286	cut clips from the ones less than 5 seconds to en-	frames from the target video and calculate	333
287	sure a strong match between video and its text. We	the average MSE (Mean Square Error) and	334
288	apply similar normalization in subsection 3.1 and	SSIM (Wang et al., 2004) between adjacent	335
289	also use SSIM and MSE between interval sampled	frames in the sample as final score.	336
290	frames to filter out the possible static videos, ensur-		
291	ing the quality in Dynamic Degree. Finally the 4K	4. Text-to-Video Alignment. We include CLIP-	337
292	real videos are scored 4 (perfect) in all aspects.	Score (Radford et al., 2021b) and X-CLIP-	338
293	We plot the rating distributions across each di-	Score (Ma et al., 2022) as metrics in this di-	339
294	mension in Figure 2. which is balanced except for	mension. CLIP-Score calculates cosine simi-	340
295	Dynamic Degree. We inspected in detail via case	larity between the feature of each frame and	341
296	study and turned out this distribution is expected.	the text prompt and then averages across all	342
297	Eventually, we get the final 37.6K examples as	frames, while X-CLIP-Score utilizes the fea-	343
298	the training split of VIDEOEVAL, and reserve 760	ture of video instead of frames.	344
299	examples as VIDEOEVAL-test for evaluation.		
300	4 Experiments	5. Factual Consistency. It is challenging to find	345
301	In this section, we describe our experiment setup,	a feature-based metric to determine whether	346
302	including baseline methods for video evaluation,	the visual content aligns with common sense.	347
303	and evaluation benchmarks for video evaluation.	Therefore, we rely on the second category of	348
304	We also discuss the training details of MANTISS-	metrics for this dimension.	349
305	CORE, and the analysis of our experiment results.		
306	4.1 Baselines	We discretized the continuous outputs of these	350
307	To compare with our evaluator model, we selected	metrics to align with our labeling scores [1, 2, 3, 4].	351
308	two categories of video quality metrics. The first	For instance, for CLIP-sim, values are converted to:	352
309	category relies on statistical or neural features for	'4' if raw output in [0.97, 1], '3' if in [0.9, 0.97),	353
310	evaluation. These metrics typically assess a single	'2' if in [0.8, 0.9) and '1' otherwise. See Table 11	354
311	video dimension such as temporal consistency,	for details.	355
312	and then yield a numerical value. The second cate-	MLLM Prompting Based Metrics To under-	356
313	gory employs advanced MLLMs to evaluate videos	stand how existing MLLMs perform on the multi-	357
314	across multiple dimensions. Extensive literature	aspect video evaluation task, we designed a prompt-	358
315	demonstrates that MLLMs not only excel in gener-	ing template in Table 9 to let them output scores	359
316	ating content on user instructions but also outper-	ranging from 1 (Bad) to 4 (Perfect) for each aspect.	360
317	form traditional metrics in evaluating AI-generated	However, some models, including Idefics2 (Lau-	361
318	content (AIGC). All baselines are listed in Table 4.	rençon et al., 2024), Fuyu (Adept AI, 2023),	362
319	Feature-Based Metrics	Kosmos-2 (Peng et al., 2023), and CogVLM (Wang	363
320	1. Visual Quality. We use two no-reference im-	et al., 2023b) and OpenFlamingo (Awadalla et al.,	364
321	age quality metrics PIQE (Venkatanath et al.,	2023), fail to give reasonable outputs. We thus ex-	365
322	2015) and BRISQUE (Mittal et al., 2012b).	clude them from the tables. MLLMs that follow the	366
323	We apply them on all frames of video and	output format like LLaVA-1.5 (Liu et al., 2023a),	367
324	take the average score across frames.	LLaVA-1.6 (Liu et al., 2024), Idefics1 (Laurençon	368
325		et al., 2023), Google's Gemini 1.5 (Reid et al.,	369
326	2. Temporal Consistency. In this dimension,	2024), and OpenAI's GPT-4o (OpenAI, 2024a).	370
327	CLIP-sim (Radford et al., 2021b) and DINO-	4.2 Evaluation Benchmarks	371
	sim (Caron et al., 2021) are computed as co-	We have included the following benchmarks to eval-	372
		uate the ability of MANTISSCORE and the above-	373
		mentioned baselines on evaluating model genera-	374
		tion results.	375

Method	Visual Quality	Temporal	Dynamic Degree	Text Alignment	Factual	Average
Random	-3.1	0.5	0.4	1.1	2.9	0.4
Feature-based automatic metrics						
PIQE	-17.7	-14.5	1.2	-3.4	-16.0	-10.1
BRISQUE	-32.4	-26.4	-4.9	-8.6	-29.1	-20.3
CLIP-sim	21.7	29.1	-34.4	2.0	26.1	8.9
DINO-sim	19.4	29.6	-37.9	2.2	24.0	7.5
SSIM-sim	33.0	30.6	-31.3	4.7	30.2	13.4
MSE-dyn	-20.3	-24.7	38.0	3.3	-23.9	-5.5
SSIM-dyn	-31.4	-29.1	31.5	-5.3	-30.0	-12.9
CLIP-Score	-10.9	-10.0	-14.7	-0.3	-0.3	-7.2
X-CLIP-Score	-3.2	-2.7	-7.3	5.9	-2.0	-1.9
MLLM Prompting						
LLaVA-1.5-7B	9.4	8.0	-2.2	11.4	15.8	8.5
LLaVA-1.6-7B	-8.0	-4.1	-5.7	1.4	0.8	-3.1
Idefics2	4.2	4.5	8.9	10.3	4.6	6.5
Gemini-1.5-Flash	24.1	5.0	20.9	21.3	32.9	20.8
Gemini-1.5-Pro	35.2	-17.2	18.2	26.7	21.6	16.9
GPT-4o	13.6	17.6	28.2	25.7	30.2	23.1
Ours						
MANTIScore (gen)	86.2	80.3	77.6	59.4	82.1	77.1
MANTIScore (reg)	84.7	81.5	68.4	59.5	84.6	75.7
Δ over Best Baseline	+51.0	+50.9	+39.6	+32.8	+51.7	+54.1

Table 4: Correlation (Spearman’s ρ) between model answer and human reference on VIDEOEVAL-test.

VIDEOEVAL-test As mentioned in section 3, we split 760 video entries from VIDEOEVAL dataset, which contains 680 annotated videos and 80 augmented videos. We take label prediction accuracy and Spearman’s ρ in each dimension as evaluation indicators. For a specific aspect in the VIDEOEVAL-test (e.g. Visual Quality), we use the predicted score from the same aspect to measure the performance for baselines and our models.

GenAI-Bench GenAI-Bench (Jiang et al., 2024b) is a benchmark designed to evaluate MLLM’s ability on preference comparison for tasks including text-to-video generation and others. The preference data is taken from GenAI-Arena from user voting. We select the video preference data in our experiments. This involves the MLLM judging which of the two provided videos is generally better, measured by pairwise accuracy. We use the averaged scores of the five aspects for MLLM prompting baselines and our models to give the preference. We compute the correlation between model-assigned preference vs. human preference as our indicator.

VBench VBench (Huang et al., 2023) is a comprehensive multi-aspect benchmark suite for video generative models, where they use a bunch of existing auto-metrics in each aspect. VBench have released a set of human pref-

erence annotations on all the aspects, comprising videos by 4 models, including ModelScope (Wang et al., 2023a), CogVideo (Hong et al., 2022), VideoCrafter1 (Chen et al., 2023a), and LaVie (Wang et al., 2023c). We select the subset from 5 aspects of VBench, like technical quality, subject consistency, and so on, to compute the preference comparison accuracy. For each aspect, we subsample 100 unique prompts in the testing. We use the averaged scores of the five aspects for MLLM prompting baselines and our models to predict the preference.

EvalCrafter EvalCrafter (Liu et al., 2023b) is a text-to-video benchmark across four dimensions: Video Quality, Temporal Consistency, Text-to-Video Alignment, and Motion Quality. We focused on the first three ones and gathered 2,541 videos by five models: Pika, Gen2, Floor33 (Floor33, 2024), ModelScope, and ZeroScope (Sterling, 2024). In EvalCrafter, human annotators rated each video on a scale of 1-5, with each scored by three raters. We calculated the average score across raters and normalized it to [0, 1]. After inference on benchmark videos, we excluded "Dynamic Degree" and "Factual Consistency" to match EvalCrafter’s dimensions. Finally, we used Spearman’s ρ in each dimension as an indicator.

Benchmark →	GenAI-Bench		VBench			
Model ↓ Sub-Aspect →	Video Preference	Technical Quality	Subject Consistency	Dyanmic Degree	Motion Smoothness	Overall Consistency
Random	37.7	44.5	42.0	37.3	40.5	44.8
Feature-based Automatic Metrics						
PIQE	34.5	60.8	44.3	71.0	45.3	53.8
BRISQUE	38.5	56.7	41.2	75.5	41.2	54.2
CLIP-sim	34.1	47.8	46.0	34.8	44.7	44.2
DINO-sim	31.4	49.5	51.2	24.7	55.5	41.7
SSIM-sim	28.4	30.7	46.2	24.5	54.2	27.2
MSE-dyn	34.2	32.8	31.7	81.7	31.2	39.2
SSIM-dyn	38.5	37.5	36.3	84.2	34.7	44.5
CLIP-Score	45.0	57.8	46.3	71.3	47.0	52.2
X-CLIP-Score	41.4	44.0	38.0	51.0	28.7	39.0
MLLM Prompting						
LLaVA-1.5-7B	49.9	42.7	42.3	63.8	41.33	8.8
LLaVA-1.6-7B	44.5	38.7	26.8	56.5	28.5	43.2
Idefics1	34.6	20.7	22.7	54.0	27.3	33.7
Gemini-1.5-Flash	67.1	52.3	49.2	64.5	45.5	49.9
Gemini-1.5-Pro	60.9	56.7	43.3	65.2	43.0	56.3
GPT-4o	52.0	59.3	49.3	46.8	42.0	60.8
Ours						
MANTISSCORE (gen)	59.0	64.2	57.7	55.5	54.3	61.5
MANTISSCORE (reg)	78.5	78.2	71.5	68.0	74.0	69.0
Δ over Best Baseline	+11.4	+17.3	+20.3	-16.2	+18.5	+8.2

Table 5: Pairwise preference accuracy on GenAI-Bench (Jiang et al., 2024b) and VBench (Huang et al., 2023). For MLLM prompting and our method, we averaged the five aspect scores defined in Table 2 as the score for each video in the comparison, where the higher one deemed the winner. The table below shows the accuracy of each method by comparing these computed scores with human annotations of "Win," "Tie," and "Lost" for the two videos.

Method	Visual	Temporal	Text Align
Random	-2.0	1.4	-0.9
EvalCraft (GPT-4V)	55.4	56.7	32.3
Feature-based Automatic Metrics			
PIQE	0.5	-3.3	-0.9
BRISQUE	6.4	-1.3	6.7
CLIP-sim	36.0	53.5	19.2
DINO-sim	30.6	50.3	15.3
SSIM-im	32.4	36.9	11.4
MSE-dyn	-15.4	-27.5	-8.1
SSIM-dyn	-32.6	-33.9	-12.6
CLIP-Score	18.7	11.5	35.0
X-CLIP-Score	12.2	3.1	24.5
MLLM Prompting			
LLaVA-1.5-7B	13.4	15.6	2.6
LLaVA-1.6-7B	12.2	8.5	18.9
Idefics1	1.5	-1.5	0.8
Gemini-1.5-Flash	34.9	-27.8	44.8
Gemini-1.5-Pro	37.8	-24.1	55.1
GPT-4o	32.9	12.5	40.7
Ours			
MANTISSCORE (gen)	20.8	51.3	10.7
MANTISSCORE (reg)	42.4	51.3	59.5
Δ over Best Baseline	-13.1	-5.4	4.4

Table 6: Spearman’s Correlation (ρ) of MANTISSCORE on EvalCrafter (Liu et al., 2023b)

4.3 Training Details

For MANTISSCORE, We use two scoring methods: generative scoring and regression scoring. Generative scoring involves training the model to output fixed text forms, from which aspect scores are extracted using regular expressions. These scores are integers corresponding to human annotation scores. In contrast, regression scoring replaces the language model head with a linear layer that outputs 5 logits representing scores for each aspect. Regression scoring is trained using MSE loss.

We select Mantis-Idefics2-8B (Jiang et al., 2024a) as the base model, which can accommodate 128 video frames at most. The learning rate is set to $1e-5$. Each model is trained for 1 epoch on 8 A100 (80G) GPUs, finishing in 6 hours.

4.4 Results

We report the Spearman correlation results on the VIDEOEVAL-test and EvalCrafter in Table 4 and Table 6, respectively. For the preference comparison on videos, we report the pairwise accuracy on the GenAI-Bench and VBench in Table 5.

Base Model	Scoring Type	VIDEOEVAL *	EvalCrafter*	GenAI-Bench	VBench*	Average
VideoLLaVA-7B	Generation	71.9	9.8	42.6	46.5	42.7
Idefics2-8B	Generation	73.9	11.3	50.7	53.9	47.5
Mantis-Idefics2-8B	Generation	<u>77.1</u>	<u>27.6</u>	59.0	58.7	55.6
Idefics2-8B	Regression	<u>73.9</u>	<u>17.4</u>	74.5	64.4	57.5
Mantis-Idefics2-8B	Regression	75.7	51.1	78.5	73.0	69.6

Table 7: Ablation study on the base model and scoring function for MANTISSCORE. "*" means that we take the average of Spearman correlation or pairwise accuracy across the multiple aspects of the benchmark. The highest numbers are bold for each benchmark, and the second are underlined.

MANTISSCORE achieves the SoTA performance

On the VIDEOEVAL-test, MANTISSCORE gets an average of 54.1 improvements on all the five aspects compared to the baseline GPT-4o. What’s more, on the EvalCrafter benchmark, MANTISSCORE (reg) has 4.4 improvements on text-to-video alignment. For pairwise preference comparison, MANTISSCORE also gets 78.5 accuracy on GenAI-Bench, surpassing the second-best Gemini-1.5-Flash by 11.4 points. on the VBench, our model archives the highest pairwise accuracy on 4 out of 5 aspects from VBench, with an average of 16.1 improvements.

Feature-based Automatic Metrics are limited

While some feature-based automatic metrics are good at a single aspect, they might fail to evaluate well on others. For example, on the VIDEOEVAL-test, the correlation scores of SSIM-dyn and MSE-dyn achieve 31.5 and 38.0 for the dynamic degree aspect, but they both get a negative correlation for others. Besides, PIQE, BRISQUE, CLIP-Score, and X-CLIP-Score get nearly all negative correlations for all 5 aspects. This proves the image quality assessment metrics cannot be easily adapted to the video quality assessment task.

4.5 Ablation Study

We conducted an ablation study on the base model selection and scoring types by training different variants on VIDEOEVAL. Results shown in Table 7.

Base model ablation To investigate the effects of changing the base model, we have trained different variants with VideoLLaVA-7B and Idefics2-8B as the base models. Since VIDEOEVAL-test, EvalCrafter, and VBench both have multiple aspects in the benchmarks, we take the average score across these aspects and report the general performance in Table 7. The results show that the Video-LLaVA-based version gets the worst performance on the four benchmarks, even if it is specifically designed for video understanding. The Idefics2-8B-based version has marginal improve-

ments compared to the VideoLLaVA. After changing to Mantis-Idefics2-8B, the scores on the four benchmarks keep improving from 47.5 to 55.6 on average. When the scoring type is regression, the mantis-based version is still better than the Idefics2-based version by 12.1 points. Therefore, we select the Mantis-based version as the final choice.

Regression scoring or generative scoring?

The primary difference between regression scoring and generative scoring is that regression scoring can give more fine-grained scores instead of just the four labels. Results on EvalCrafter, GenAI-Bench, and VBench all indicate that using regression scoring can consistently improve the Spearman correlation or the pairwise comparison accuracy. For example, on GenAI-Bench, MANTISSCORE (reg) achieves 78.5 accuracy, which is higher than the 59.0 of the MANTISSCORE (gen). The results are similar for the other benchmarks. We thus conclude that regression scoring with more fine-grained scores is a better choice.

5 Conclusion

In this paper, we introduce MANTISSCORE, which is trained on our meticulously curated dataset VIDEOEVAL for video evaluation. We hired 20 expert raters to annotate the 37.6K videos generated from 11 popular text-to-video generative models across 5 key aspects, Visual Quality, Temporal Consistency, Dynamic Degree, Text-to-Video Alignment and Factual Consistency. Our IAA match ratio gets more than 60%. We test the performance of MANTISSCORE using Spearman correlation on VIDEOEVAL-test and EvalCrafter, and using pairwise comparison accuracy on GenAI-Bench and VBench. The results show that MANTISSCORE consistently gets the best performance, surpassing the powerful baseline GPT-4o and Gemini 1.5 Flash/Pro by a large margin. Our work highlights the importance of using MLLM for video evaluation due to its rich world knowledge and the high-quality rating dataset across multiple aspects.

535

References

536

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

541

Adept AI. 2023. Fuyu-8B: A Multimodal Architecture for AI Agents. <https://www.adept.ai/blog/fuyu-8b>.

544

Jie An, Songyang Zhang, Harry Yang, Sonal Gupta, Jia-Bin Huang, Jiebo Luo, and Xi Yin. 2023. Latent-shift: Latent diffusion with temporal shift for efficient text-to-video generation. *arXiv preprint arXiv:2304.08477*.

549

Anas Awadalla, Irena Gao, Josh Gardner, Jack Hes-sel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jernia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. 2023. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*.

557

Hritik Bansal, Zongyu Lin, Tianyi Xie, Zeshun Zong, Michal Yarom, Yonatan Bitton, Chenfanfu Jiang, Yizhou Sun, Kai-Wei Chang, and Aditya Grover. 2024. Videophy: Evaluating physical commonsense for video generation.

562

Omer Bar-Tal, Hila Chefer, Omer Tov, Charles Her-mann, Roni Paiss, Shiran Zada, Ariel Ephrat, Jun-hwa Hur, Yuanzhen Li, Tomer Michaeli, et al. 2024. Lumiere: A space-time diffusion model for video generation. *arXiv preprint arXiv:2401.12945*.

567

Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendeleevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. 2023a. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*.

573

Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. 2023b. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575.

580

Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. 2021. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660.

586

Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, et al. 2023a. Videocrafter1: Open diffusion models for high-quality video generation. *arXiv preprint arXiv:2310.19512*.

591

Haoxin Chen, Yong Zhang, Xiaodong Cun, Menghan Xia, Xintao Wang, Chao Weng, and Ying Shan. 2024a. Videocrafter2: Overcoming data limitations for high-quality video diffusion models. *arXiv preprint arXiv:2401.09047*.

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647	Roberto Henschel, Levon Khachatryan, Daniil Hayrapetyan, Hayk Poghosyan, Vahram Tadevosyan, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. 2024. Streamingt2v: Consistent, dynamic, and extendable long video generation from text. <i>arXiv preprint arXiv:2403.14773</i> .	701
648		702
649		703
650		704
651		
652		
653	Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. <i>Advances in neural information processing systems</i> , 33:6840–6851.	705
654		706
655		707
656		708
657	Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. 2022. Cogvideo: Large-scale pre-training for text-to-video generation via transformers. <i>arXiv preprint arXiv:2205.15868</i> .	709
658		710
659		711
660		
661	Berthold KP Horn and Brian G Schunck. 1981. Determining optical flow. <i>Artificial intelligence</i> , 17(1-3):185–203.	712
662		713
663		714
664	Zhewei Huang, Tianyuan Zhang, Wen Heng, Boxin Shi, and Shuchang Zhou. 2022. Real-time intermediate flow estimation for video frame interpolation. <i>Preprint</i> , arXiv:2011.06294.	715
665		716
666		717
667		
668	Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, Yao-hui Wang, Xinyuan Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. 2023. Vbench: Comprehensive benchmark suite for video generative models. <i>ArXiv</i> , abs/2311.17982.	718
669		719
670		720
671		
672		
673		
674		
675	Dongfu Jiang, Xuan He, Huaye Zeng, Cong Wei, Max Ku, Qian Liu, and Wenhui Chen. 2024a. Mantis: Interleaved multi-image instruction tuning. <i>arXiv preprint arXiv:2405.01483</i> .	721
676		722
677		723
678		724
679	Dongfu Jiang, Max Ku, Tianle Li, Yuansheng Ni, Shizhuo Sun, Rongqi Fan, and Wenhui Chen. 2024b. Genai arena: An open evaluation platform for generative models. <i>arXiv preprint arXiv:2406.04485</i> .	725
680		
681		
682		
683	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. <i>arXiv preprint arXiv:2303.13439</i> .	726
684		727
685		728
686		729
687		730
688		
689	Tengchuan Kou, Xiaohong Liu, Zicheng Zhang, Chunyi Li, Haoning Wu, Xiongkuo Min, Guangtao Zhai, and Ning Liu. 2024a. Subjective-aligned dataset and metric for text-to-video quality assessment. <i>ArXiv</i> , abs/2403.11956.	731
690		732
691		733
692		734
693		735
694	Tengchuan Kou, Xiaohong Liu, Zicheng Zhang, Chunyi Li, Haoning Wu, Xiongkuo Min, Guangtao Zhai, and Ning Liu. 2024b. Subjective-aligned dataset and metric for text-to-video quality assessment. <i>arXiv preprint arXiv:2403.11956</i> .	736
695		737
696		738
697		739
698		
699	Klaus Krippendorff. 2011. Computing krippendorff’s alpha-reliability.	740
700		
	Max Ku, Dongfu Jiang, Cong Wei, Xiang Yue, and Wenhui Chen. 2023. Viescore: Towards explainable metrics for conditional image synthesis evaluation. <i>Preprint</i> , arXiv:2312.14867.	701
		702
		703
		704
	Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu Cord, and Victor Sanh. 2023. Obelics: An open web-scale filtered dataset of interleaved image-text documents. <i>Preprint</i> , arXiv:2306.16527.	705
		706
		707
		708
		709
		710
		711
	Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. 2024. What matters when building vision-language models? <i>ArXiv</i> , abs/2405.02246.	712
		713
		714
	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024. Llava-next: Improved reasoning, ocr, and world knowledge.	715
		716
		717
	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023a. Visual instruction tuning. <i>ArXiv</i> , abs/2304.08485.	718
		719
		720
	Yaofang Liu, Xiaodong Cun, Xuebo Liu, Xintao Wang, Yong Zhang, Haoxin Chen, Yang Liu, Tiejong Zeng, Raymond Chan, and Ying Shan. 2023b. Evalcrafter: Benchmarking and evaluating large video generation models.	721
		722
		723
		724
		725
	Yuanxin Liu, Lei Li, Shuhuai Ren, Rundong Gao, Shicheng Li, Sishuo Chen, Xu Sun, and Lu Hou. 2023c. Fetv: A benchmark for fine-grained evaluation of open-domain text-to-video generation. <i>Preprint</i> , arXiv:2311.01813.	726
		727
		728
		729
		730
	Yiwei Ma, Guohai Xu, Xiaoshuai Sun, Ming Yan, Ji Zhang, and Rongrong Ji. 2022. X-clip: End-to-end multi-grained contrastive learning for video-text retrieval. In <i>Proceedings of the 30th ACM International Conference on Multimedia</i> , pages 638–647.	731
		732
		733
		734
		735
	Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. 2012a. No-reference image quality assessment in the spatial domain. <i>IEEE Transactions on image processing</i> , 21(12):4695–4708.	736
		737
		738
		739
	Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. 2012b. No-reference image quality assessment in the spatial domain. <i>IEEE Transactions on Image Processing</i> , 21(12):4695–4708.	740
		741
		742
		743
	John Mullan, Duncan Crawbuck, and Aakash Sastry. 2023. Hotshot-XL.	744
		745
	OpenAI. 2024a. Hello gpt-4o. https://openai.com/index/hello-gpt-4o/ . Accessed: 2024-06-15.	746
		747
	OpenAI. 2024b. Video generation models as world simulators.	748
		749
	Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. Kosmos-2: Grounding multimodal large language models to the world. <i>ArXiv</i> , abs/2306.14824.	750
		751
		752
		753

754	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021a. Learning transferable visual models from natural language supervision. In <i>International Conference on Machine Learning</i> .	810
755		811
756		812
757		
758		813
759		814
760		815
761	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021b. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR.	816
762		817
763		
764		818
765		819
766		820
767		821
768	Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	822
769		823
770		824
771		825
772		826
773		
774	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 10684–10695.	827
775		828
776		829
777		830
778		831
779		832
780	Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. 2016. Improved techniques for training gans. <i>Advances in neural information processing systems</i> , 29.	833
781		834
782		835
783		836
784	Spencer Sterling. 2024. Zeroscope v2 .	837
785		
786	Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. 2018. Towards accurate generative models of video: A new metric & challenges. <i>arXiv preprint arXiv:1812.01717</i> .	
787		
788		
789		
790	Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. 2019. FVD: A new metric for video generation .	
791		
792		
793		
794	N. Venkatanath, D. Praneeth, Maruthi Chandrasekhar Bh, Sumohana Channappayya, and Swarup Medasani. 2015. Blind image quality evaluation using perception based features . <i>2015 21st National Conference on Communications, NCC 2015</i> .	
795		
796		
797		
798		
799		
800	Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. 2023a. Modelscope text-to-video technical report . <i>ArXiv</i> , abs/2308.06571.	
801		
802		
803		
804	Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. 2023b. Cogvlm: Visual expert for pretrained language models . <i>ArXiv</i> , abs/2311.03079.	
805		
806		
807		
808		
809		
	Wenhao Wang and Yi Yang. 2024. Vidprom: A million-scale real prompt-gallery dataset for text-to-video diffusion models . <i>ArXiv</i> , abs/2403.06098.	
	Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yanan He, Jiashuo Yu, Peiqing Yang, et al. 2023c. Lovie: High-quality video generation with cascaded latent diffusion models. <i>arXiv preprint arXiv:2309.15103</i> .	
	Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. <i>IEEE transactions on image processing</i> , 13(4):600–612.	
	Haoning Wu, Chaofeng Chen, Jingwen Hou, Liang Liao, Annan Wang, Wenxiu Sun, Qiong Yan, and Weisi Lin. 2022. Fast-vqa: Efficient end-to-end video quality assessment with fragment sampling. In <i>European Conference on Computer Vision</i> , pages 538–554.	
	Haoning Wu, Erli Zhang, Liang Liao, Chaofeng Chen, Jingwen Hou, Annan Wang, Wenxiu Sun, Qiong Yan, and Weisi Lin. 2023. Exploring video quality assessment on user generated contents from aesthetic and technical perspectives. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 20144–20154.	
	Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. 2022. Magicvideo: Efficient video generation with latent diffusion models. <i>arXiv preprint arXiv:2211.11018</i> .	

838	A Ethical Statement		884
839	This work fully complies with the ACL Ethics Policy. We declare that there are no ethical issues in this paper, to the best of our knowledge.		885
840			886
841			887
842	B Risks and Limitation		888
843	Although we have designed systematic pipelines to recruit expert raters and annotate the video evaluation scores, we still find out that some annotations contain errors and may harm the overall quality of the dataset. Our IAA score computation is only based on a small number of trial examples and, thus might not represent the actual IAA of the whole annotations.		889
844			890
845			891
846			892
847			893
848			894
849			895
850			896
851	Besides, while MANTISSCORE is proven to be able to effectively give reasonable scores on our defined five aspects, it can still sometimes output wrong scores that do not match our expectations. We admit this drawback and list that as one of our future works.		897
852			898
853			899
854			900
855			901
856			902
857	C Dataset Licence		903
858	We have used VidProM (Wang and Yang, 2024) to collect the prompts used for video generation, whose usage LICENSE is CC BY-NC 4.0 license. For other evaluation datasets, We did not find license for EvalCrafter (Liu et al., 2023b) human annotations. GenAI-Bench (Jiang et al., 2024b) is under MIT licence, and VBench (Huang et al., 2023) is under Apache 2.0 license. We are thus able to utilize these datasets in our experiments.		904
859			905
860			906
861			907
862			908
863			909
864			910
865			911
866			912
867			913
868			914
869			915
870	D Annotator Management		916
871	During the annotation, we have recruited 20 expert raters, where 14 of them are undergraduate or graduate students, who will become one of the authors of our paper, and the rest of them are assured to be paid with decent salary.		917
872			918
873			919
874			920
875			921
876	E Video Format Normalizing Details		922
877	To mitigate difference of videos format from different generative models, we normalize the frame rate of all the generated videos to 8 fps (frames per second). Specifically, for high frame rate model Pika and AnimateDiffusion (Guo et al., 2023), we use uniform down-sampling to normalize Pika from 24 fps to 8fps, and Animate-		923
878			924
879			925
880			926
881			927
882			928
883			929
			930
			931
			932
			933
			934
			935
			936
			937
			938
			939
			940
			941
			942
			943
			944
			945
			946
			947
			948
			949
			950
			951
			952
			953
			954
			955
			956
			957
			958
			959
			960
			961
			962
			963
			964
			965
			966
			967
			968
			969
			970
			971
			972
			973
			974
			975
			976
			977
			978
			979
			980
			981
			982
			983
			984
			985
			986
			987
			988
			989
			990
			991
			992
			993
			994
			995
			996
			997
			998
			999
			1000

Videos Gallery -- See examples in each sub-score

1. Visual quality

Expected Case:

- (1) The video looks clear and normal on its appearance.
- (2) The features like Brightness, Contrast, Color, etc, are appropriate and stable.

Error point:

- (a) local obvious unclear or blurry,
- (b) too low resolution,
- (c) some speckles or black patches,
- (d) appearance of video is skewed and distorted,
- (e) unstable optical property, such as brightness, contrast, saturation, exposure etc,
- (f) flickering color of main objects and background

Note:

**Some videos have watermark, we can ignore that.

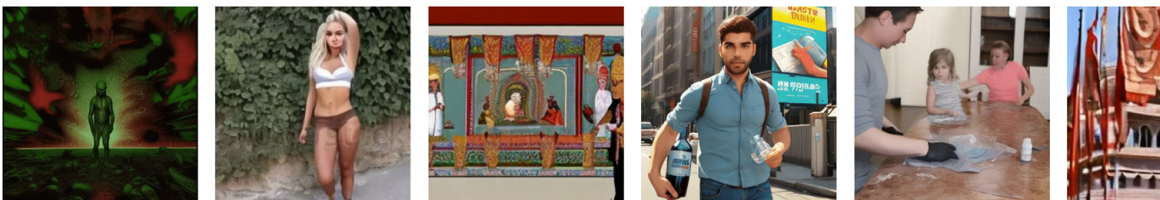
Annotator Login

Username:

Visual Quality - Good



Visual Quality - Avg



Visual Quality - Bad



Figure 3: Welcome Page of our video annotating website, with definition, checklist for error points and diverse video examples.

Report Problem

Log out

(1) You can submit answers for several times, we'll take the final version.

(2) **scoring standards**

Good: The video matches the "Expected case" very well, there is no "error point" and video quality is not affected

Average: There is one not too severe "error point", which have affected the video quality to some extent.

Bad: The video is very inconsistent with "Expected case". There are one or more obviously severe "error points" that have greatly affected the video quality

1. Visual quality

Expected Case:

(1) The video looks clear and normal on its appearance.

(2) The features like Brightness, Contrast, Color, etc, are appropriate and stable.

Error point:

(a) local obvious unclear or blurry,

(b) too low resolution,

(c) some speckles or black patches,

(d) appearance of video is skewed and distorted,

(e) unstable optical property, such as brightness, contrast,

Least index of incomplete video : 1

Current progress: 26/1000

Turn to:


1. visual quality Good Average Bad

2. temporal consistency Good Average Bad

3. dynamic degree Good Average Bad

4. text alignment Good Average Bad

5. fact consistency Good Average Bad



Input Prompt:
image realistic, tarot cards in hands of young woman, prime plane, close up image, very realistic

Answered.

(*3: Good, *2: Average, *1: Bad
*-1: reported as problematic.)

visual/optical quality	2	temporal consistency	3
dynamic degree	3	text-to-video alignment	3
factual consistency	1		

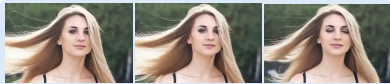
Figure 4: Working page of our video annotating website

I Case study of VIDEOEVAL

We showcase the annotations examples in Figure 5. The first example depicts a clear video of a woman with her hair moving, thus scoring 3 in all 5 aspects. The second example shows a distorted video, thus scoring 1 across all the aspects except the dynamic degree. We further analyzed the correlations between the designed aspects in Figure 6. We found that visual quality achieves a high correlation of 0.6 with temporal consistency, while dynamic degree has a very low correlation with all other aspects.

Video Annotation

Text: The wind gently blows, and her hair moves. Nothing else moves.



Visual Quality : 3

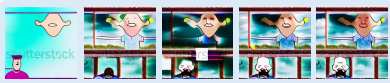
Temporal Consistency : 3

Dynamic Degree : 3

Text-to-Video Alignment : 3

Factual Consistency : 3

Text: A transformation of a happy, energetic cartoon character gradually transitioning to a state of burnout and exhaustion.



Visual Quality : 1

Temporal Consistency : 1

Dynamic Degree : 3

Text-to-Video Alignment : 1

Factual Consistency : 1

Figure 5: Example of annotations. Each video has a text description and is rated for the 5 aspects.

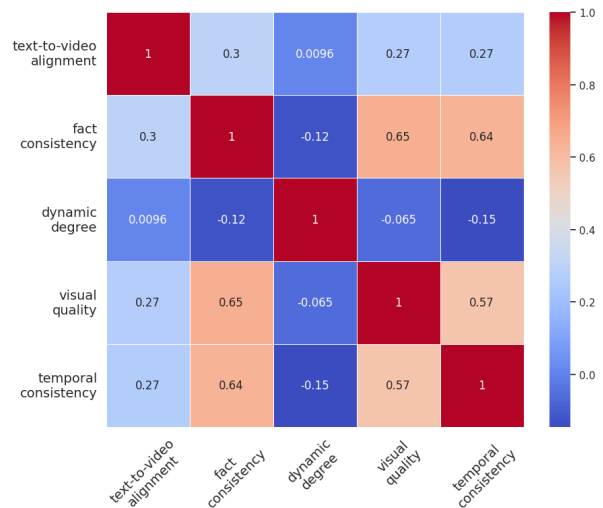


Figure 6: Correlation study on the evaluation aspects.

931

932

933

934

935

936

937

938

939

940

941

<p>Evaluation Aspect</p> <p>Visual Quality</p>	<p>Detailed Description for Annotation</p> <p>Expected Case:(1) The video looks clear and normal on its appearance. (2) The features like Brightness, Contrast, Color, etc, are appropriate and stable.</p> <p>Error point:(a) local obvious unclear or blurry, (b) too low resolution, (c) some speckles or black patches, (d) appearance of video is skewed and distorted, (e) unstable optical property, such as brightness, contrast, saturation, exposure etc, (f) flickering color of main objects and background</p> <p>Note:Some videos have watermark, we can ignore that.</p>
<p>Temporal Consistency</p>	<p>Expected Case: (1) The main objects, main characters and overall appearance are consistent across the video. (2) The appearance of video as well as the movements of humans and objects are smooth and natural.</p> <p>Error points: (a) The person or object suddenly disappears or appears (b) The type or class of objects has obvious changes (c) There is an obvious switch in the screen shot (d) the appearance of video or movements in it is laggy and un-smooth, (e) local deformation or dislocation of human or objects due to the motion (for large scale deformation, the video should also be rated as bad in "1. visual quality"),</p> <p>Note: For a video almost static or with small dynamic degree, as long as it does not have error points, then it should be scored as good.</p>
<p>Dynamic Degree</p>	<p>Expected Case: (1) The video is obviously not static, the people or objects or the video screen is dynamic. (2) The video can be easily distinguished from a static image.</p> <p>Note: You are supposed to focus on only dynamic degree, regardless of the visual quality and video content</p>
<p>Text-to-Video Alignment</p>	<p>Expected Case: The characters, objects, motions, events etc. that are mentioned in text input prompts all exist reasonably.</p> <p>Error points: (a) The people and objects in prompt do not appear in video (b) The actions and events in prompt do not appear in video (c) The number, size, shape, color, state, movement and other attributes of the objects in the video do not match the prompt (d) Text mentioned in prompt is not displayed correctly in the video, such as "a placard saying 'No Smoking'" but "No Smoking" is not spelled correctly in the video (e) The video format (such as width, height, screen ratio, duration) does not match the format in prompt.</p>
<p>Factual Consistency</p>	<p>Expected Case: (1) Overall appearance and motion are consistent with our common-sense, physical principles, moral standards, etc.</p> <p>Error points: (a) static ones: Content in video goes against common sense in life, such as lighting a torch in the water, standing in the rain but not getting wet, etc. (b) static ones: The size, color, shape and other basic properties of objects violate scientific principles (c) dynamic ones: The overall movement of people or objects violates common-sense and laws of physics, such as spontaneous upward movement against gravity, abnormal water flow, etc. (d) dynamic ones: Partial movements of people or objects violate common-sense and laws of physics, such as the movement of hands or legs is anti-joint, etc.</p> <p>Notes: Relation with '5. text-to-video alignment': Some text prompts express fictional and unrealistic content, for example, "a dog plays the guitar in the sky" or "an astronaut rides a horse in space". In this case, regardless of the veracity of the text prompt, you should only consider whether the other content in the video makes sense.</p>

Table 8: Expected cases and error cases for each aspect that annotators can see during the annotation.

Suppose you are an expert in judging and evaluating the quality of AI-generated videos, please watch the following frames of a given video and see the text prompt for generating the video, then give scores from 5 different dimensions:

- (1) visual quality: the quality of the video in terms of clearness, resolution, brightness, and color
- (2) temporal consistency, the consistency of objects or humans in video
- (3) dynamic degree, the degree of dynamic changes
- (4) text-to-video alignment, the alignment between the text prompt and the video content
- (5) factual consistency, the consistency of the video content with the common-sense and factual knowledge

For each dimension, output a number from [1,2,3,4], in which '1' means 'Bad', '2' means 'Average', '3' means 'Good', '4' means 'Real' or 'Perfect' (the video is like a real video)

Here is an output example:

visual quality: 4
temporal consistency: 4
dynamic degree: 3
text-to-video alignment: 1
factual consistency: 2

For this video, the text prompt is "{text_prompt}", all the frames of video are as follows:

Table 9: Prompting template in generation format used for MANTISSCORE training and the MLLM prompting baselines

Suppose you are an expert in judging and evaluating the quality of AI-generated videos, please watch the following frames of a given video and see the text prompt for generating the video, then give scores from 5 different dimensions:

- (1) visual quality: the quality of the video in terms of clearness, resolution, brightness, and color
- (2) temporal consistency, the consistency of objects or humans in video
- (3) dynamic degree, the degree of dynamic changes
- (4) text-to-video alignment, the alignment between the text prompt and the video content
- (5) factual consistency, the consistency of the video content with the common-sense and factual knowledge

For each dimension, output a float number from 1.0 to 4.0, higher the number is, better the video performs in that dimension, the lowest 1.0 means Bad, the highest 4.0 means Perfect/Real (the video is like a real video)

Here is an output example:

visual quality: 2.24
temporal consistency: 3.89
dynamic degree: 3.17
text-to-video alignment: 1.86
factual consistency: 2.16

For this video, the text prompt is "{text_prompt}", all the frames of video are as follows:

Table 10: Prompting template used for the MLLM prompting baseline and MANTISSCORE training

Dimension	Metric	1 (Bad)	2 (Avg)	3 (Good)	4 (Perfect)
Visual Quality	PIQE↓	[50, ∞)	[30,50)	[15,30)	[0,15]
	BRISQUE↓	[50,∞)	[30,50)	[10,30)	[0,10]
Temporal Consistency	CLIP-sim↑	[0,0.80)	[0.80,0.90)	[0.90,0.97)	[0.97,1]
	DINO-sim↑	[0,0.75)	[0.75,0.85)	[0.85,0.95)	[0.95,1]
	SSIM-sim↑	[0,0.6)	[0.6,0.75)	[0.75,0.9)	[0.9,1]
Dynamic Degree	MSE-dyn↑	[0,100]	[100,1000)	[1000,3000)	[3000,∞)
	SSIM-dyn↓	[0.9,1]	[0.7,0.9)	[0.5,0.7)	[0,0.5)
Text-to-Video Alignment	CLIP-Score↑	[0.2,0.27)	[0.27,0.31)	[0.31,0.35)	[0.35,0.4]
	X-CLIP-Score↑	[0,0.15)	[0.15,0.23)	[0.23,0.30)	[0.30,1]

Table 11: Discretization rules for featured-based baselines.