

# AIGVE-MACS: Unified Multi-Aspect Commenting and Scoring Model for AI-Generated Video Evaluation

Anonymous CVPR submission

Paper ID

## Abstract

001 *The rapid advancement of AI-generated video models has*  
002 *created a pressing need for robust and interpretable eval-*  
003 *uation frameworks. Existing metrics are limited to produc-*  
004 *ing numerical scores without explanatory comments, result-*  
005 *ing in low interpretability and human evaluation alignment.*  
006 *To address those challenges, we introduce AIGVE-MACS, a*  
007 *unified model for AI-Generated Video Evaluation(AIGVE),*  
008 *which can provide not only numerical scores but also multi-*  
009 *aspect language comment feedbacks in evaluating these*  
010 *generated videos. Central to our approach is AIGVE-*  
011 *BENCH 2, a large-scale benchmark comprising 2,500 AI-*  
012 *generated videos and 22,500 human-annotated detailed*  
013 *comments and numerical scores across nine critical evalua-*  
014 *tion aspects. Leveraging AIGVE-BENCH 2, AIGVE-MACS*  
015 *incorporates recent Vision-Language Models with a novel*  
016 *token-wise weighted loss and a dynamic frame sampling*  
017 *strategy to better align with human evaluators. Compre-*  
018 *hensive experiments across supervised and zero-shot bench-*  
019 *marks demonstrate that AIGVE-MACS achieves state-of-*  
020 *the-art performance in both scoring correlation and com-*  
021 *ment quality, significantly outperforming prior baselines in-*  
022 *cluding GPT-4o and VideoScore. In addition, we further*  
023 *showcase a multi-agent refinement framework where feed-*  
024 *back from AIGVE-MACS drives iterative improvements in*  
025 *video generation, leading to 53.5% quality enhancement.*  
026 *This work establishes a new paradigm for comprehensive,*  
027 *human-aligned evaluation of AI-generated videos.*

## 028 1. Introduction

029 With the rapid advancement of video generation models  
030 such as Sora [35], HunyuanVideo [17], and Mochi-1 [38],  
031 AI-generated videos are becoming increasingly photoreal-  
032 istic and temporally coherent, significantly narrowing the  
033 gap between synthetic and real-world visual content. Fol-  
034 lowing provided textual instruction prompts, AI-generated  
035 videos are gaining widespread adoption across domains

such as entertainment, advertising, education, and virtual 036  
reality, offering cost-effective solutions and enabling more 037  
personalized content experiences for both creators and con- 038  
sumers [41, 51]. 039

Despite significant progress, AI-generated videos still 040  
face persistent challenges such as limited spatial resolu- 041  
tion, object distortions, and misalignment with user instruc- 042  
tions [11, 22, 42]. These issues highlight the growing im- 043  
portance of rigorous evaluation, which plays a critical role 044  
in guiding the development and refinement of video gener- 045  
ation models. However, research on evaluation methodolo- 046  
gies has lagged behind the rapid advances in video gener- 047  
ation itself, leaving a gap in systematic, interpretable, and 048  
comprehensive evaluation techniques [28]. 049

Existing evaluation methods for AI-generated videos ex- 050  
hibit several key limitations. First, they often rely on out- 051  
dated numerical metrics such as FVD [39] and IS [4], which 052  
are insufficient for capturing the nuanced qualities of mod- 053  
ern AI-generated content. Effective AI-Generated Video 054  
Evaluation (AIGVE) must account not only for intrinsic 055  
video quality but also for alignment with user-provided in- 056  
structions [28, 46]. Therefore, evaluating video quality with 057  
one single scalar metric is inherently limiting. While recent 058  
work has made progress by decomposing quality into as- 059  
pects such as technical quality, motion dynamics, and text- 060  
to-video alignment [3, 7, 11], many approaches still rely on 061  
aggregating legacy metrics not originally designed for this 062  
task [14, 29]. These metrics often overlap in what they mea- 063  
sure or fail to address key dimensions altogether, resulting 064  
in evaluation pipelines that are fragmented, difficult to in- 065  
terpret, and lacking comprehensive coverage. 066

Second, although recent studies have explored the use 067  
of unified Vision Language Models (VLMs) for evaluat- 068  
ing AI-generated videos, they often fail to fully leverage 069  
the generative strengths of these models. This shortcoming 070  
stems largely from the difficulty of reliably prompting or 071  
finetuning VLMs to produce accurate, human-aligned eval- 072  
uative outputs. Rather than generating natural language 073  
assessments directly, many existing approaches either at- 074  
tach a scoring head to predict numerical values from the 075

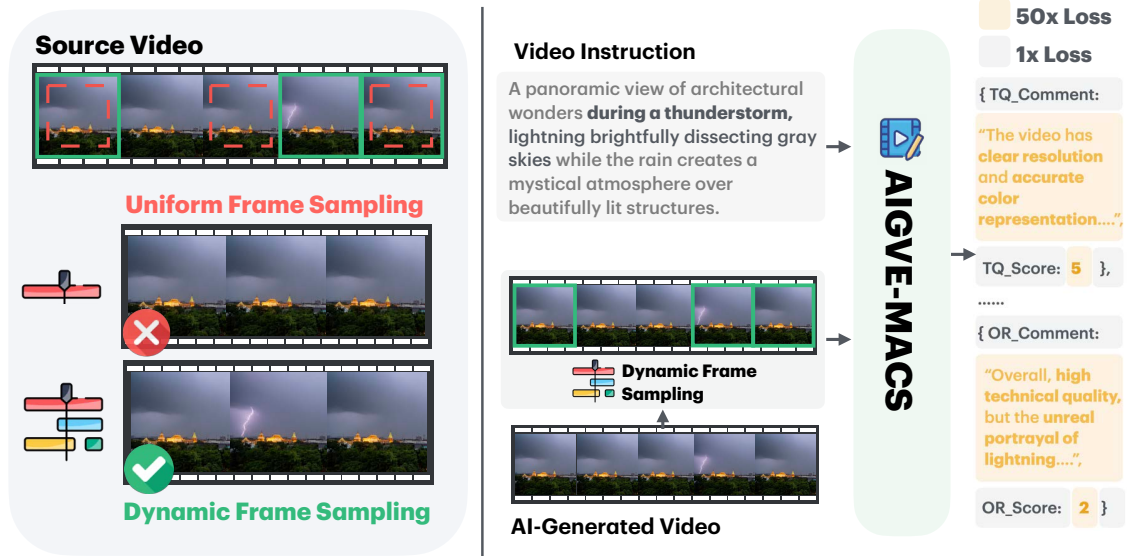


Figure 1. AIGVE-MACS Pipeline. The left side of the diagram illustrates our Dynamic Frame Sampling strategy, which selects key moments based on content variation to better capture temporal dynamics. The right side highlights the Token-Wise Weighted Loss, which emphasizes score and comment tokens to improve alignment with human evaluations. TQ and OR refer to Technical Quality and Overall, respectively. Tokens in the yellow box are assigned a loss weight of 50, while those in the gray box are assigned a loss weight of 1.

076 VLM’s hidden states, thereby bypassing its language generation capabilities [10], or reformulate evaluation as a series of binary yes/no questions [7, 49]. These strategies reduce the VLM to a classification tool, limiting its ability to provide nuanced, context-aware evaluations. Furthermore, the question-based approach relies heavily on handcrafted prompts, which struggle to capture the complexity and diversity of real-world video quality dimensions. Besides, due to the constrained context length of current VLMs, it is impossible to input the full token sequence of an entire video. Instead, existing methods typically sample a fixed number of frames uniformly across the video, which can overlook brief but critical dynamic changes.

089 Third, due to the absence of datasets that pair aspect-wise scores with explanatory comments, existing video evaluation approaches primarily focus on producing numerical scores [8, 10, 19, 32] while neglecting the generation of human-like comments that offer richer, more actionable feedback. However, we argue that such comments are actually essential for understanding the specific strengths and weaknesses of AI-generated videos. By providing detailed qualitative insights, explanatory comments can guide model improvement more effectively than scalar scores alone, which will also be demonstrated with show-cases to be provided in this paper.

101 To address the aforementioned challenges, we first introduce and release the AIGVE-BENCH 2, a large-scale, human-annotated benchmark comprising 500 diverse prompts, 2,500 videos generated by five state-of-the-art

105 video generation models, and 22,500 high-quality human score and comment annotations spanning nine carefully designed evaluation aspects. The prompts, videos, scores, and comments in AIGVE-BENCH 2 are diligently constructed and rigorously validated to ensure comprehensive coverage of diverse video generation scenarios, making the dataset both robust and representative of real-world applications.

112 Building on this benchmark, we propose to finetune Qwen2.5-VL-7B [2] and present AIGVE-MACS, a unified evaluation model that provides both numerical score and natural language comment feedback evaluations across the nine critical aspects as covered in the AIGVE-BENCH 2. The overall pipeline of the proposed AIGVE-MACS is also illustrated in Figure 1. By introducing a dynamic frame sampling strategy as well as a weighted loss on both comment and score tokens, AIGVE-MACS achieves accurate numerical predictions while fully leveraging the generative capabilities of VLMs to generate language comment to justify those scores. This joint modeling approach enhances the transparency, interpretability, and practical utility of automated video evaluation. Extensive experiments and ablation studies show that AIGVE-MACS achieves state-of-the-art performance in both supervised and zero-shot settings on AIGVE-BENCH 2 and multiple evaluation benchmarks, significantly outperforming existing methods in both numerical scoring and comment generation.

131 What’s more, to further demonstrate the practical utility of our model, we introduce a multi-agent iterative refinement framework that leverages the generated scores and

134	comments to progressively improve video quality. Experi-	185
135	mental results indicate that the quality of generated videos	186
136	can be enhanced by 53.5% through this iterative process.	187
137	This framework underscores the real-world applicability of	188
138	AIGVE-MACS, illustrating how continuous feedback loops	189
139	can drive meaningful improvements in AI-generated con-	190
140	tent. To the best of our knowledge, this is the first work	191
141	to jointly produce aspect-wise scores and natural language	192
142	comments for AI-generated videos, providing a more holistic	
143	evaluation framework that aligns closely with human	
144	judgments.	
145	In summary, our contributions are as follows:	
146	• We construct AIGVE-BENCH 2, the first large-scale	193
147	AIGVE benchmark covering nine evaluation aspects,	
148	with 22,500 human-annotated numerical scores and ex-	194
149	planatory comments for 2,500 videos generated by sev-	195
150	eral SOTA video generation models.	196
151	• We propose AIGVE-MACS, a unified AIGVE model	197
152	that jointly predicts aspect-wise scores and comments.	198
153	Trained with dynamic frame sampling and weighted	199
154	losses, AIGVE-MACS achieves state-of-the-art perfor-	200
155	mance on AIGVE-BENCH 2 and other benchmarks in	201
156	both supervised and zero-shot settings, aligning well with	202
157	human judgments.	203
158	• We explore practical applications of our model through a	204
159	multi-agent iterative refinement framework that uses gen-	205
160	erated scores and comments to enhance video quality.	
161	<b>2. Related Works</b>	
162	<b>2.1. AI-Generated Video Evaluation</b>	
163	The field of AI-Generated Video Evaluation (AIGVE) is	206
164	in its early stages and includes even more challenges than	207
165	video quality assessment. Previous research proposes that	208
166	evaluating AI-generated videos requires alignment with	209
167	both human perception and instructions to ensure high-	210
168	quality video generation that meets creators' intentions and	211
169	viewers' expectations [28]. Alignment with human percep-	212
170	tion focuses on evaluating video quality through traditional	213
171	metrics like resolution and clarity, while ensuring consis-	214
172	tency with physical world properties such as realistic tex-	215
173	tures and adherence to physical laws. Meanwhile, align-	216
174	ment with human instructions emphasizes how well videos	217
175	mirror the detailed scenarios, actions, and narratives speci-	218
176	fied in text descriptions, ensuring the generated content ful-	219
177	fills creators' creative and communicative objectives. These	220
178	aspects have evolved separately, with perception alignment	221
179	progressing from statistical approaches [4, 39] to advanced	222
180	architectures like DOVER [45] for aesthetic assessment and	223
181	BVQI [44] for CLIP-based evaluation, while instruction	224
182	alignment developed methods like TIFA [13] and CLIP-	225
183	Score [12] for evaluating text-visual consistency.	
184	Recent comprehensive frameworks attempt to unify	
	these aspects. VBench [14] evaluates across 16 dimensions,	226
	and EvalCrafter [29] assesses four key aspects including vi-	227
	sual quality and text-video alignment. VIDEOSCORE [10]	228
	leverages large language models for holistic assessment.	229
	However, these methods either rely on collections of tradi-	230
	tional metrics or provide single scores without interpretable	231
	feedback, highlighting the need for more sophisticated eval-	232
	uation approaches.	233
	<b>3. AIGVE-BENCH 2</b>	234
	Our AIGVE-BENCH 2 extends the previous benchmark	
	dataset AIGVE-BENCH [46], which contains multi-aspect	
	numerical scores for AI-generated videos across nine eval-	
	uation dimensions. We introduce a rigorous comment	
	process pipeline to enrich AIGVE-BENCH with human-	
	written, aspect-specific justifications. This pipeline ensures	
	that each textual explanation is directly aligned with the cor-	
	responding multi-aspect numerical evaluation scores, offer-	
	ing a more comprehensive and interpretable understanding	
	of video quality. The remainder of this section will first	
	present the key characteristics of AIGVE-BENCH and	
	then detail our comment processing pipeline.	
	<b>3.1. Previous Benchmark: AIGVE-BENCH</b>	
	AIGVE-BENCH [46] is a comprehensive AIGVE bench-	
	mark dataset contains 2,500 AI-generated videos created	
	from 500 diverse prompts, covering a broad range of real-	
	world scenes and interactions. Each prompt falls into one	
	of two categories: <i>global view</i> , which focuses on envi-	
	ronmental and wide-angle scenes with natural dynamics	
	such as weather changes, camera motion, and <i>close shot</i> ,	
	which captures fine-grained interactions between humans,	
	animals, plants, or objects. Videos are generated by five	
	state-of-the-art models: Pyramid [15], CogVideoX [47],	
	Genmo [38], Hunyuan [17], and Sora [35], with varying res-	
	olutions and frame rates, and are standardized to a uniform	
	duration of 5 seconds to ensure sufficient temporal depth	
	and comparability.	
	Each video is evaluated by expert raters across nine dis-	
	tinct evaluation aspects: Technical Quality, Dynamics, Con-	
	sistency, Physics, Element Presence, Element Quality, Ac-	
	tion/Interaction Presence, Action/Interaction Quality, and	
	Overall.	
	<b>3.2. Comment Processing Pipeline</b>	
	Although justification comments are collected in parallel	
	with the scoring process, they present several challenges	
	that limit their direct usability. First, since the comments are	
	written by multiple annotators, they exhibit inconsistencies	
	in writing style, terminology, and level of detail, introducing	
	noise and increasing the complexity of model training. Sec-	
	ond, comments may not fully align with scores. As shown	
	in Figure 2, minor issues may be emphasized without noting	

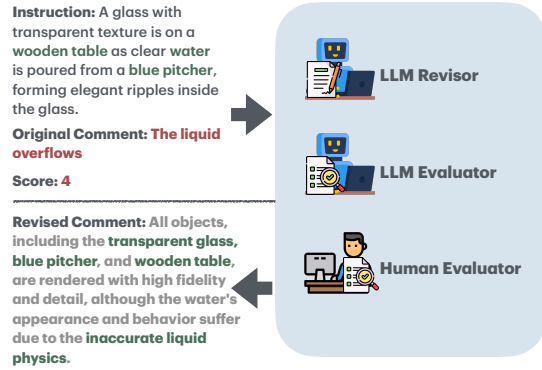


Figure 2. Overview of the comment processing pipeline for AIGVE-BENCH 2. The pipeline includes a revisor, validator, and human evaluator to ensure high-quality, consistent, and objective comments.

235 strengths, leading to insufficient justification for high scores  
236 and reducing their value for training or explanation.

237 To address these challenges, we design a robust comment  
238 processing pipeline that revises the original comments to  
239 ensure consistency, objectivity, and alignment with the cor-  
240 responding scores. As shown in Figure 2, we first employ  
241 a Large Language Model (LLM) as a revisor to refine and  
242 extend the original comments, conditioned on the video in-  
243 structions and associated scores. The LLM is prompted to  
244 produce comments with a consistent and objective writing  
245 style, providing a comprehensive explanation that covers  
246 both the strengths and weaknesses relevant to the evalua-  
247 tion aspect.

248 Next, the revised comment is passed to a second LLM,  
249 which we refer to as the evaluator. This model checks  
250 whether the revised comment introduces any content that  
251 are not presented in the original comment or instruction,  
252 mitigating the risk of hallucination. Finally, a human eval-  
253 uator manually reviews and, if necessary, rewrites the revised  
254 comments to ensure clarity, factual accuracy, and faithful re-  
255 flection of the corresponding scores. Through this meticu-  
256 lous process, we extend each aspect-level score in AIGVE-  
257 BENCH with a high-quality justification comment, resulting  
258 in a total of 22,500 comments with an average length of 267  
259 words.

## 260 4. AIGVE-MACS

261 To build the AIGVE-MACS framework, we propose to fine-  
262 tune the VLMs with our AIGVE-BENCH 2, enabling VLMs  
263 to generate human-aligned evaluation scores and explanat-  
264 ory comments. Specifically, given an AI-generated video,  
265 we provide the model with the video instruction and video  
266 content. The VLM is trained to generate a JSON-formatted  
267 output containing both the aspect-specific score and the cor-  
268 responding natural language comment.

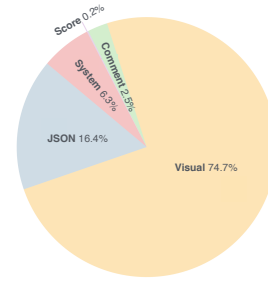


Figure 3. Averaged ratios of different types of input tokens in the input sequence. The five types of tokens are: (1) system prompt tokens, (2) visual tokens, (3) JSON structure tokens, (4) comment tokens, and (5) numerical score tokens.

Formally, given an AI-generated video  $v$ , the video’s in- 269  
struction  $i$ , the VLM is trained to generate a structured out- 270  
put: 271

$$y = \text{VLM}(i, v|\theta), \quad (1) \quad 272$$

where  $y = \{C, S\}$  denotes the JSON-formatted out- 273  
put containing aspect-specific evaluation comments  $C$  and 274  
scores  $S$ ,  $\theta$  is the learnable parameter of the VLM. 275

276 However, as discussed in Section 1, existing methods are  
277 facing two major challenges. First, they struggle to effec-  
278 tively leverage the native language generation capabilities  
279 of VLMs to produce human-aligned evaluation scores and  
280 comments. Second, the commonly adopted uniform frame  
281 sampling strategy may overlook brief yet critical dynamics,  
282 leading to incorrect evaluation. To address these issues, we  
283 first introduce a token-wise weighted loss in Section 4.1,  
284 which encourages the model to focus more on generating  
285 informative comments and accurate scores by highlighting  
286 these tokens during training. Additionally, in Section 4.2,  
287 we propose a dynamic frame sampling strategy to better  
288 capture key moments in the video.

### 289 4.1. Token-wise Weighted Loss

290 To better understand the challenge of leveraging language  
291 generation ability, we analyze the input tokens used by  
292 VLMs in our task. Specifically, there are five types of in-  
293 put tokens: (1) system prompt tokens that describe the task,  
294 (2) visual tokens that encode the video content, (3) JSON  
295 structure tokens that maintain the output format, (4) com-  
296 ment tokens that contain user comments, and (5) score to-  
297 kens that indicate the evaluation scores. Figure 3 illustrates  
298 the averaged ratios of each type of input token in the input.  
299 We observe that the system prompt tokens, visual tokens,  
300 and JSON structure tokens, which are not directly related  
301 to the task, occupy a significant portion of the input, while  
302 the task-related comment tokens and score tokens are rela-  
303 tively sparse. This unbalanced distribution of input tokens  
304 will lead to the VLMs focusing more on the system prompt,

305 visual content, and output format, while neglecting the user  
306 comments and scores.

307 To address this issue, we propose a simple yet effective  
308 token-wise weighted loss that encourages the model to fo-  
309 cus on accurately generate comment and score tokens. Dur-  
310 ing tokenization, we identify comment and score tokens and  
311 assign them a higher weight  $\alpha$  in the loss function. The  
312 weighted loss is defined as follows:

$$313 \quad \mathcal{L} = - \sum_{t=1}^T w_t \cdot \log p(y_t | y_{<t}, i, v), \quad (2)$$

314 where  $T$  is the total number of tokens in the output se-  
315 quence,  $y_t$  is the target token at position  $t$ ,  $x$  is the input  
316 sequence, and  $p(y_t | y_{<t}, i, v)$  is the predicted probability  
317 of token  $y_t$ . The token weight  $w_t$  is defined as:

$$318 \quad w_t = \begin{cases} \alpha, & \text{if } y_t \text{ is a comment or score token,} \\ 1, & \text{otherwise.} \end{cases} \quad (3)$$

319 Besides, inspired by reasoning schema, we put the com-  
320 ment content prior to the numerical score in training data to  
321 encourage the model to build stronger connections between  
322 the comment and the score to enhance the comment-score  
323 alignment. The process is shown on the right part of Fig-  
324 ure 1.

## 325 4.2. Dynamic Frame Sampling

326 To avoid missing the key moments in the video, we pro-  
327 pose a dynamic frame sampling strategy that selects frames  
328 based on content variation instead of uniform sampling.  
329 As shown at the bottom-left of Figure 1, this approach al-  
330 lows the evaluation model to focus on the most informative  
331 frames, leading to more accurate and aspect-aware assess-  
332 ments.

333 To select representative frames based on content  
334 changes, we compute a frame-wise difference score as fol-  
335 lows:

$$336 \quad \Delta(f^t, f^{t-1}) = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W \mathbb{I}(|f_{i,j}^t - f_{i,j}^{t-1}| > 0) \quad (4)$$

337 where  $f^t \in \mathbb{R}^{H \times W}$  denotes the grayscale frame with  
338 width  $W$  and height  $H$  at time  $t$ , and  $\mathbb{I}(\cdot)$  is the indicator  
339 function. A frame  $f^t$  is selected if it satisfies both:

$$340 \quad \Delta(f^t, f^{t-1}) > \theta \quad \text{and} \quad t - t_{\text{last}} \geq \gamma, \quad (5)$$

341 where  $\theta$  is a content change threshold and  $\gamma$  is the minimum  
342 frame gap to prevent redundant selection. If the number of  
343 selected frames exceeds the target number  $N$ , we uniformly  
344 subsample them. If no frames meet the criteria, we fallback  
345 to uniform sampling.

## 5. Experiment 346

### 5.1. Benchmark Datasets and Evaluation Metrics 347

348 We evaluate AIGVE-MACS on four benchmarks: a su-  
349 pervised dataset, AIGVE-BENCH 2-TEST, and three zero-  
350 shot datasets—VIDEOFEEBACK-Test [10], GenAI-  
351 Bench [23], and VBench [14], using the same test samples  
352 as VideoScore [10]. For AIGVE-BENCH 2-TEST, we re-  
353 port Spearman’s  $\rho$  for score prediction and use ROUGE-1,  
354 ROUGE-L [25], UniEval-Fact [52], BERTScore [50], and  
355 G-Eval [30] for comment evaluation. VIDEOFEEBACK-  
356 Test uses Spearman’s  $\rho$  over five aspects; GenAI-Bench and  
357 VBench reports pairwise accuracy inferred from predicted  
358 scores.

### 5.2. Baseline Methods 359

360 To benchmark the performance of AIGVE-MACS, we com-  
361 pare it against a comprehensive suite of existing video eval-  
362 uation methods, organized into three main categories:

**Feature-based Metrics** These methods evaluate videos 363  
364 by extracting handcrafted or pretrained model features.  
365 This category includes traditional quality metrics such as  
366 PIQE [40] and BRISQUE [33], perceptual similarity-based  
367 scores like SSIM-sim [43], SSIM-dyn [43], MSE(Mean  
368 Square Error)-dyn, and DINO-sim [5], as well as text-video  
369 alignment metrics such as CLIP-sim [18], CLIP-temp [29],  
370 CLIP-Score [12], X-CLIP-Score [31], BLIP-sim [24], and  
371 PickScore [16].

**Modeling-based Metrics** This group includes supervised 372  
373 deep learning models trained specifically to predict video  
374 quality. We include LightVQA+ [53], GSTVQA [6], and  
375 SimpleVQA [9], which are adapted for multi-aspect video  
376 assessment.

**VLM-based Metrics** These approaches leverage VLMs 377  
378 to generate scores or conduct evaluations via either prompt-  
379 ing or finetuning. Baselines include foundation mod-  
380 els such as Qwen2.5-VL [1], VideoLLaMA3-7B [48],  
381 GPT-4.1 [36], GPT-4o [34], LLaVA-1.5-7B [26], LLaVA-  
382 1.6-7B [27], Gemini-1.5-Flash [37], Gemini-1.5-Pro [37],  
383 Idedics1 [20], and Idedics2 [21]. In addition, we consider  
384 VLM-based evaluation frameworks like VideoScore [10],  
385 and TIFA [13].

### 5.3. Implementation Details 386

387 We implement AIGVE-MACS based on the Qwen2.5-VL-  
388 7B [2] architecture, initializing from the official check-  
389 point<sup>1</sup>. During finetuning, the vision encoder is kept frozen,

<sup>1</sup><https://huggingface.co/Qwen/Qwen2.5-VL-7B-Instruct>

Method	Technical	Dynamic	Consistency	Physics	Element Pre	Element Qu	Act Pre	Act Qu	Overall	AVG
Random	-6.71	4.28	-3.55	-1.62	-3.03	1.42	-5.16	1.44	-3.19	-1.79
<b>Feature-based Metrics</b>										
CLIP-sim	9.12	6.54	5.79	-4.45	19.34	7.93	23.85	8.25	21.58	10.88
BLIP-sim	12.02	10.92	12.24	3.69	22.80	9.42	17.02	10.34	19.09	13.06
CLIP-temp	16.46	4.26	<u>27.69</u>	<u>23.51</u>	11.07	25.12	-2.62	20.20	16.63	15.81
PickScore	22.48	5.90	11.36	6.55	26.25	17.69	20.37	13.24	24.28	16.46
<b>Modeling-based Metrics</b>										
LightVQA+	-3.68	-7.58	-8.08	-10.73	-0.16	-6.33	4.77	-7.40	-4.70	-4.88
GSTVQA	17.92	13.40	15.97	1.23	-2.65	15.91	9.81	9.68	20.71	11.33
SimpleVQA	24.50	11.50	16.58	0.28	2.41	18.30	7.76	3.58	21.22	11.79
<b>VLM-based Metrics</b>										
Qwen2.5-VL	8.77	4.00	1.24	-6.01	9.19	10.19	18.74	0.72	9.59	6.27
VideoLLaMA3	15.94	<u>19.44</u>	11.70	13.21	-3.13	12.27	13.61	-0.69	11.58	10.44
GPT-4.1	<u>36.49</u>	5.81	26.68	19.87	<u>27.22</u>	28.77	32.75	20.22	29.98	25.31
GPT-4o	34.71	7.05	18.12	20.28	23.10	<u>30.47</u>	<u>36.57</u>	<u>31.58</u>	<u>38.57</u>	<u>26.72</u>
Videoscore	-9.50	-8.20	-0.20	20.10	9.70	-7.50	-3.10	-0.60	-7.30	-0.73
VideoPhy	0.10	4.00	-1.40	-5.60	0.80	-1.30	11.90	2.70	9.70	2.32
TIFA	17.81	8.83	9.62	3.85	16.67	12.41	17.90	5.87	17.78	12.30
AIGVE-MACS	<b>40.60</b>	<b>57.31</b>	<b>61.49</b>	<b>64.36</b>	<b>40.32</b>	<b>40.81</b>	<b>44.31</b>	<b>60.71</b>	<b>59.88</b>	<b>52.20</b>

Table 1. Scoring Correlation Evaluation Result on AIGVE-BENCH 2. The underlined score represents the best zero-shot model. Element Pre and Element Qu represents Element Presence and Quality. Act Pre and Qu represents Action Presence and Quality.

Method	GenAI-Bench (Acc)	VBench (Acc)						VideoFeedback ( $\rho$ )
		Technical Quality	Subject Consistency	Dynamic Degree	Motion Smoothness	Overall Consistency	Average	
Random	37.7	44.5	42.0	37.3	40.5	44.8	41.82	0.4
<b>Feature-based Metrics</b>								
PIQE	34.5	60.8	44.3	71.0	45.3	53.8	55.04	-10.1
BRISQUE	38.5	56.7	41.2	75.5	41.2	54.2	53.76	-20.3
CLIP-sim	34.1	47.8	46.0	34.8	44.7	44.2	43.50	8.9
DINO-sim	31.4	49.5	51.2	24.7	55.5	41.7	44.52	7.5
SSIM-sim	28.4	30.7	46.2	24.5	54.2	27.2	36.56	13.4
MSE-dyn	34.2	32.8	31.7	81.7	31.2	39.2	43.32	-5.5
SSIM-dyn	38.5	37.5	36.3	<u>84.2</u>	34.7	44.5	47.44	-12.9
CLIP-Score	45.0	57.8	46.3	71.3	47.0	52.2	54.92	-7.2
X-CLIP-Score	41.4	44.0	38.0	51.0	28.7	39.0	40.14	-1.9
<b>VLM-based Metrics</b>								
LLaVA-1.5	49.9	42.7	42.3	63.8	41.3	8.8	39.78	8.5
LLaVA-1.6	44.5	38.7	26.8	56.5	28.5	43.2	38.74	-3.1
Idefics1/2	34.6	20.7	22.7	54.0	27.3	33.7	31.68	6.5
Gemini-1.5-Flash	67.1	52.3	49.2	64.5	45.5	49.9	52.28	20.8
Gemini-1.5-Pro	60.9	56.7	43.3	65.2	43.0	56.3	52.90	16.9
GPT-4o	52.0	59.3	49.3	46.8	42.0	60.8	51.64	23.0
VideoScore	59.0	64.2	57.7	55.5	54.3	<u>61.5</u>	58.64	77.1
AIGVE-MACS (Ours)	<b>68.5</b>	<b>65.5</b>	<b>59.0</b>	<b>74.5</b>	<b>55.8</b>	<b>55.6</b>	<b>62.08</b>	<b>24.6</b>

Table 2. Evaluation results on GenAI-Bench, VBench, and VideoFeedback. GenAI and VBench report pairwise accuracy (%), while VideoFeedback uses Spearman’s  $\rho$ . The VBench section spans metrics from Technical Quality to Average. Underlined scores indicate best zero-shot performance.

390 while all remaining parameters are updated. We set the  
 391 weight for score and comment to 50. Training is conducted  
 392 on two NVIDIA A6000 GPUs using a learning rate of 1e-  
 393 5. We train the model for 3 epochs and apply early stop-  
 394 ping based on validation loss. Optimization is performed  
 395 using AdamW, with a linear learning rate scheduler and  
 396 5% warmup steps. The full training process takes approxi-  
 397 mately 24 hours.

## 5.4. Results

AIGVE-MACS achieves state-of-the-art alignment with human judgments, significantly outperforming all existing baselines in both score prediction and comment generation. Table 1 and Table 3 summarize AIGVE-MACS’s performance on AIGVE-BENCH 2 compared to a broad range of evaluation baselines. AIGVE-MACS achieves the highest correlation across all nine evaluation aspects, with an average Spearman’s  $\rho$  of 52.20%, which nearly doubling the performance of the strongest VLM baseline, GPT-4o, and achieving up to 61% improvement over its pretrained

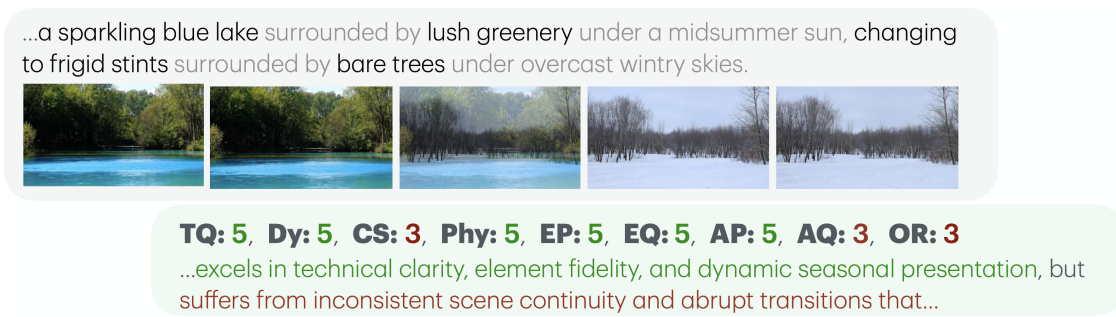


Figure 4. Case Study of AIGVE-MACS. This example showcases AIGVE-MACS’s ability to produce human-aligned evaluations and distinguish fine-grained quality differences in AI-generated videos.

	R1	RL	BS	UF	GE
GPT-4o	18.30	15.86	74.90	40.84	2.10
GPT-4.1	15.80	12.94	73.99	43.99	2.10
QWen2.5-VL	17.95	15.31	74.31	42.32	2.37
VideoLLama	19.99	17.67	75.35	40.21	2.18
AIGVE-MACS	<b>49.50</b>	<b>38.00</b>	<b>85.87</b>	<b>57.04</b>	<b>3.42</b>

Table 3. Evaluation Results for Comments. R1, RL, BS, UF, and GE represent ROUGE-1, ROUGE-L, BERTcore, UniEval-Fact, and G-Eval respectively.

409 backbone, Qwen2.5-VL-7B [2].

410 In terms of comment generation, as shown in Table 3,  
 411 AIGVE-MACS also consistently outperforms all baselines  
 412 across standard automatic metrics. These results highlight a  
 413 fundamental limitation of current VLMs, which struggle to  
 414 generate reliable, human-aligned evaluations. In contrast,  
 415 AIGVE-MACS effectively learns to generate both accurate  
 416 scores and rich, faithful commentary that aligns closely with  
 417 human judgment.

418 **AIGVE-MACS establishes new state-of-the-art perfor-**  
 419 **mance across all three benchmarks, demonstrating**  
 420 **strong zero-shot generalization capabilities for diverse**  
 421 **video evaluation tasks.** Table 2 reports zero-shot evalua-  
 422 tion results on VideoFeedback [11], GenAI-Bench [23], and  
 423 VBench [14].

424 On GenAI-Bench, AIGVE-MACS achieves the high-  
 425 est video preference accuracy, outperforming strong com-  
 426 mercial VLMs such as Gemini-1.5-Pro and GPT-4o. On  
 427 VBench, AIGVE-MACS surpasses VideoScore by 4.24%,  
 428 while using only one-tenth of the training data.

429 Notably, on the VideoFeedback dataset, AIGVE-MACS  
 430 also achieves significant improvements over all open-source  
 431 and commercial baselines, emerging as the best zero-shot  
 432 model. By comparing the result in Table 1 and Table 2,  
 433 while AIGVE-MACS generalizes well to VideoFeedback  
 434 without training on it, VideoScore, which is trained on  
 435 VideoFeedback, fails to generalize to our AIGVE-BENCH

2-TEST. This contrast highlights AIGVE-MACS’s more ro-  
 436 bust and transferable evaluation capabilities that are not tied  
 437 to the idiosyncrasies of any single dataset.  
 438

The superior supervised and zero-shot performance of  
 439 AIGVE-MACS validates the effectiveness of our finetuning  
 440 strategy, which equips VLMs with the ability to evaluate  
 441 AI-generated videos using both structured scores and natu-  
 442 ral language comments. It also underscores the high quality  
 443 and broad coverage of AIGVE-BENCH 2, enabling the train-  
 444 ing of generalizable video evaluators that can seamlessly  
 445 transfer to other benchmarks.  
 446

## 6. Multi-Agent Iterative Refinement Frame- 447 work 448

Our AIGVE-MACS opens new possibilities for leveraging  
 449 evaluation scores and comments to guide and optimize  
 450 the video generation process. To initiate this exploration,  
 451 we propose a multi-agent iterative refinement framework  
 452 designed to progressively enhance video quality through  
 453 feedback-driven revisions.  
 454

As shown in Figure 5, our framework consists of three  
 455 main components: a Video Generator, an Instruction Revi-  
 456 sor, and AIGVE-MACS. In each iteration, the Video Gen-  
 457 erator produces a video based on the current user instruc-  
 458 tion, and then AIGVE-MACS evaluates the generated video  
 459 and provides feedback in the form of scores and comments.  
 460 The Instruction Revisor refines the instruction based on the  
 461 feedback, aiming to clarify or adjust the requirements for  
 462 the next iteration. This process continues until the video  
 463 reaches a satisfactory quality level or a maximum of 4 iter-  
 464 ations.  
 465

To validate the effectiveness of our multi-agent frame-  
 466 work, we use HunyuanVideo [17] as the Video Generator,  
 467 GPT-4.1 [36] as the Instruction Revisor, and AIGVE-MACS  
 468 as the evaluator. We conduct experiments on a subset of  
 469 50 low-quality videos (overall score < 3) sampled from  
 470 AIGVE-BENCH 2-TEST. The iteration continues until the  
 471 overall score exceeds 4 or the iteration limit of 4 is reached.  
 472

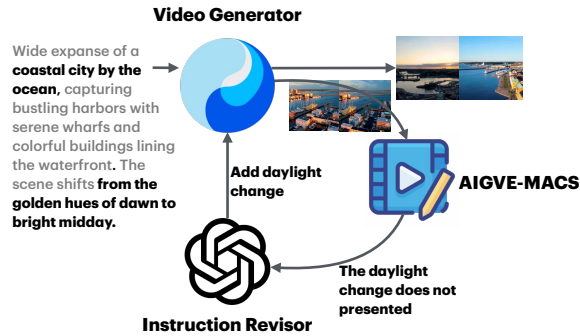


Figure 5. The pipeline of the multi-agent iterative refinement framework.

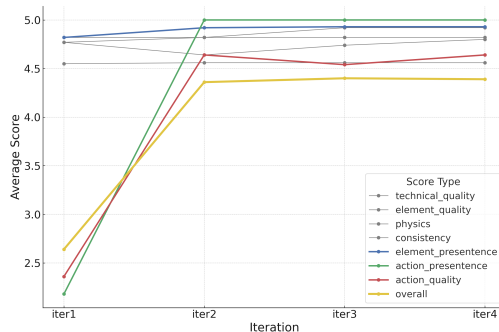


Figure 6. The result of the multi-agent iterative refinement framework.

Figure 6 shows the average scores across refinement iterations. Notably, aspects related to *overall* score and instruction alignment such as element presence, action presence, and action quality, demonstrate significant improvements up to 53%, particularly within the first two iterations.

In contrast, visual aspects like technical quality, element quality, and physics show only minor variation, as they are primarily constrained by the fixed video generator. This suggests that while our framework effectively improves instruction-driven dimensions, visual quality remains bounded by generator capacity.

## 7. Conclusion

We propose AIGVE-MACS, a unified and interpretable evaluation framework that jointly predicts aspect-wise scores and explanatory comments for AI-generated videos. Trained on our human-annotated AIGVE-BENCH 2 dataset, the model achieves state-of-the-art performance in both supervised and zero-shot settings, outperforming strong baselines in score correlation and comment quality across multiple benchmarks. Beyond evaluation, we demonstrate the utility of AIGVE-SCORE in a multi-agent refinement framework, showing its ability to drive meaningful improvements in video generation quality through feedback-

driven instruction updates. Our work highlights the value of structured evaluation signals and establishes AIGVE-MACS as a robust and generalizable tool for aligning video generation with human preferences.

## References

- [1] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-VL: A Versatile Vision-Language Model for Understanding, Localization, Text Reading, and Beyond, 2023. arXiv:2308.12966 [cs]. 5
- [2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 2, 5, 7
- [3] Hritik Bansal, Zongyu Lin, Tianyi Xie, Zeshun Zong, Michal Yarom, Yonatan Bitton, Chenfanfu Jiang, Yizhou Sun, Kai-Wei Chang, and Aditya Grover. Videophy: Evaluating physical commonsense for video generation. *arXiv preprint arXiv:2406.03520*, 2024. 1
- [4] Shane Barratt and Rishi Sharma. A note on the inception score. *arXiv preprint arXiv:1801.01973*, 2018. 1, 3
- [5] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 5
- [6] Baoliang Chen, Lingyu Zhu, Guo Li, Fangbo Lu, Hongfei Fan, and Shiqi Wang. Learning generalized spatial-temporal deep feature representation for no-reference video quality assessment. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(4):1903–1916, 2022. 5
- [7] Rui Chen, Lei Sun, Jing Tang, Geng Li, and Xiangxiang Chu. Finger: Content aware fine-grained evaluation with reasoning for ai-generated videos. *arXiv preprint arXiv:2504.10358*, 2025. 1, 2
- [8] Zijian Chen, Wei Sun, Yuan Tian, Jun Jia, Zicheng Zhang, Jiarui Wang, Ru Huang, Xiongkuo Min, Guangtao Zhai, and Wenjun Zhang. Gaia: Rethinking action quality assessment for ai-generated videos, 2024. 2
- [9] Xianfu Cheng, Wei Zhang, Shiwei Zhang, Jian Yang, Xiangyuan Guan, Xianjie Wu, Xiang Li, Ge Zhang, Jiaheng Liu, Yuying Mai, Yutao Zeng, Zhoufutu Wen, Ke Jin, Baorui Wang, Weixiao Zhou, Yunhong Lu, Tongliang Li, Wenhao Huang, and Zhoujun Li. Simplevqa: Multimodal factuality evaluation for multimodal large language models, 2025. 5
- [10] Xuan He, Dongfu Jiang, Ge Zhang, Max Ku, Achint Soni, Sherman Siu, Haonan Chen, Abhranil Chandra, Ziyang Jiang, Aaran Arulraj, Kai Wang, Quy Duc Do, Yuansheng Ni, Bohan Lyu, Yaswanth Narsupalli, Rongqi Fan, Zhiheng Lyu, Bill Yuchen Lin, and Wenhao Chen. VideoScore: Building automatic metrics to simulate fine-grained human feedback for video generation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 2105–2123, Miami, Florida, USA, 2024. Association for Computational Linguistics. 2, 3, 5, 1

- 551 [11] Xuan He, Dongfu Jiang, Ge Zhang, Max Ku, Achint Soni,  
552 Sherman Siu, Haonan Chen, Abhranil Chandra, Ziyang Jiang,  
553 Aaran Arulraj, Kai Wang, Quy Duc Do, Yuansheng Ni, Bo-  
554 han Lyu, Yaswanth Narsupalli, Rongqi Fan, Zhiheng Lyu,  
555 Yuchen Lin, and Wenhu Chen. Videoscore: Building auto-  
556 matic metrics to simulate fine-grained human feedback for  
557 video generation. *ArXiv*, abs/2406.15252, 2024. 1, 7
- 558 [12] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras,  
559 and Yejin Choi. CLIPScore: A reference-free evaluation  
560 metric for image captioning. In *Proceedings of the 2021*  
561 *Conference on Empirical Methods in Natural Language Pro-*  
562 *cessing*, pages 7514–7528, Online and Punta Cana, Domini-  
563 can Republic, 2021. Association for Computational Linguis-  
564 tics. 3, 5
- 565 [13] Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Os-  
566 tendorf, Ranjay Krishna, and Noah A Smith. Tifa: Accurate  
567 and interpretable text-to-image faithfulness evaluation with  
568 question answering, 2023. 3, 5
- 569 [14] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si,  
570 Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin,  
571 Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin  
572 Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. VBench: Com-  
573 prehensive Benchmark Suite for Video Generative Models,  
574 2023. *arXiv:2311.17982 [cs]*. 1, 3, 5, 7
- 575 [15] Yang Jin, Zhicheng Sun, Ningyuan Li, Kun Xu, Kun Xu,  
576 Hao Jiang, Nan Zhuang, Quzhe Huang, Yang Song, Yadong  
577 Mu, and Zhouchen Lin. Pyramidal flow matching for effi-  
578 cient video generative modeling. In *Proceedings of the In-*  
579 *ternational Conference on Learning Representations (ICLR)*,  
580 2025. 3
- 581 [16] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Ma-  
582 tiana, Joe Penna, and Omer Levy. Pick-a-pic: An open  
583 dataset of user preferences for text-to-image generation.  
584 2023. 5
- 585 [17] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo  
586 Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jian-  
587 wei Zhang, Kathrina Wu, Qin Lin, Junkun Yuan, Yanxin  
588 Long, Aladdin Wang, Andong Wang, Changlin Li, DuoJun  
589 Huang, Fang Yang, Hao Tan, Hongmei Wang, Jacob Song,  
590 Jiawang Bai, Jianbing Wu, Jinbao Xue, Joey Wang, Kai  
591 Wang, Mengyang Liu, Pengyu Li, Shuai Li, Weiyan Wang,  
592 Wenqing Yu, Xincheng Deng, Yang Li, Yi Chen, Yutao Cui,  
593 Yuanbo Peng, Zhentao Yu, Zhiyu He, Zhiyong Xu, Zixiang  
594 Zhou, Zunnan Xu, Yangyu Tao, Qinglin Lu, Songtao Liu,  
595 Dax Zhou, Hongfa Wang, Yong Yang, Di Wang, Yuhong Liu,  
596 Jie Jiang, and Caesar Zhong. Hunyuanvideo: A systematic  
597 framework for large video generative models. *arXiv preprint*  
598 *arXiv:2412.03603*, 2024. 1, 3, 7
- 599 [18] Tengchuan Kou, Xiaohong Liu, Zicheng Zhang, Chunyi Li,  
600 Haoning Wu, Xiongkuo Min, Guangtao Zhai, and Ning Liu.  
601 Subjective-aligned dataset and metric for text-to-video qual-  
602 ity assessment. In *Proceedings of the 32nd ACM Interna-*  
603 *tional Conference on Multimedia*, pages 7793–7802, 2024.  
604 5
- 605 [19] Tengchuan Kou, Xiaohong Liu, Zicheng Zhang, Chunyi Li,  
606 Haoning Wu, Xiongkuo Min, Guangtao Zhai, and Ning Liu.  
607 Subjective-Aligned Dataset and Metric for Text-to-Video  
608 Quality Assessment, 2024. *arXiv:2403.11956 [cs]*. 2
- [20] Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bek-  
man, Amanpreet Singh, Anton Lozhkov, Thomas Wang,  
Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela,  
Matthieu Cord, and Victor Sanh. Obelics: An open web-  
scale filtered dataset of interleaved image-text documents,  
2023. 5
- [21] Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor  
Sanh. What matters when building vision-language models?,  
2024. 5
- [22] Daeun Lee, Jaehong Yoon, Jaemin Cho, and Mohit Bansal.  
Videorepair: Improving text-to-video generation via mis-  
alignment evaluation and localized refinement. *arXiv*  
*preprint arXiv:2411.15115*, 2024. 1
- [23] Baiqi Li, Zhiqiu Lin, Deepak Pathak, Jiayao Li, Yixin Fei,  
Kewen Wu, Tiffany Ling, Xide Xia, Pengchuan Zhang, Gra-  
ham Neubig, and Deva Ramanan. Genai-bench: Evaluat-  
ing and improving compositional text-to-visual generation,  
2024. 5, 7, 1
- [24] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi.  
Blip: Bootstrapping language-image pre-training for unified  
vision-language understanding and generation, 2022. 5
- [25] Chin-Yew Lin. ROUGE: A package for automatic evaluation  
of summaries. In *Text Summarization Branches Out*, pages  
74–81, Barcelona, Spain, 2004. Association for Computa-  
tional Linguistics. 5, 1
- [26] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee.  
Improved baselines with visual instruction tuning, 2024. 5
- [27] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan  
Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Im-  
proved reasoning, ocr, and world knowledge, 2024. 5
- [28] Xiao Liu, Xinhao Xiang, Zizhong Li, Yongheng Wang,  
Zhuoheng Li, Zhuosheng Liu, Weidi Zhang, Weiqi Ye, and  
Jiawei Zhang. A survey of ai-generated video evaluation.  
*arXiv preprint arXiv:2410.19884*, 2024. 1, 3
- [29] Yaofang Liu, Xiaodong Cun, Xuebo Liu, Xintao Wang,  
Yong Zhang, Haoxin Chen, Yang Liu, Tiejong Zeng, Ray-  
mond Chan, and Ying Shan. EvalCrafter: Benchmarking  
and Evaluating Large Video Generation Models, 2023.  
*arXiv:2310.11440 [cs]*. 1, 3, 5
- [30] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen  
Xu, and Chenguang Zhu. G-eval: Nlg evaluation using gpt-4  
with better human alignment, 2023. 5, 1
- [31] Yiwei Ma, Guohai Xu, Xiaoshuai Sun, Ming Yan, Ji Zhang,  
and Rongrong Ji. X-clip: End-to-end multi-grained con-  
trastive learning for video-text retrieval. In *Proceedings of*  
*the 30th ACM international conference on multimedia*, pages  
638–647, 2022. 5
- [32] Yibo Miao, Yifan Zhu, Yinpeng Dong, Lijia Yu, Jun Zhu,  
and Xiao-Shan Gao. T2vsafetybench: Evaluating the safety  
of text-to-video generative models, 2024. 2
- [33] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad  
Bovik. No-reference image quality assessment in the spa-  
tial domain. *IEEE Transactions on Image Processing*, 21  
(12):4695–4708, 2012. 5
- [34] OpenAI. Gpt-4o. <https://openai.com/index/hello-gpt-4o/>, 2024. 5
- [35] OpenAI. Sora. <https://openai.com/index/sora/>, 2024. 1, 3

- 667 [36] OpenAI. Gpt-4.1. <https://openai.com/index/gpt-4-1/>, 2025. 5, 7
- 668
- 669 [37] Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context, 2024. 5
- 670
- 671 [38] Genmo Team. Mochi 1. <https://github.com/genmoai/models>, 2024. 1, 3
- 672
- 673 [39] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric and challenges, 2019. 1, 3
- 674
- 675 [40] Narasimhan Venkatanath, D Praneeth, Maruthi Chandrasekhar Bh, Sumohana S Channappayya, and Swarup S Medasani. Blind image quality evaluation using perception based features. In *2015 twenty first national conference on communications (NCC)*, pages 1–6. IEEE, 2015. 5
- 680
- 681
- 682 [41] Shreyas Verma and John Leddo. Comparing the effectiveness between human-generated videos and ai-generated videos on learning. 1
- 683
- 684
- 685 [42] Wenjing Wang, Huan Yang, Zixi Tuo, Huiguo He, Junchen Zhu, Jianlong Fu, and Jiaying Liu. Swap attention in spatiotemporal diffusions for text-to-video generation. *International Journal of Computer Vision*, pages 1–19, 2025. 1
- 686
- 687
- 688 [43] Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4): 600–612, 2004. 5
- 689
- 690
- 691
- 692 [44] Haoning Wu, Liang Liao, Annan Wang, Chaofeng Chen, Jingwen Hou, Wenxiu Sun, Qiong Yan, and Weisi Lin. Towards robust text-prompted semantic criterion for in-the-wild video quality assessment, 2023. 3
- 693
- 694
- 695 [45] Haoning Wu, Erli Zhang, Liang Liao, Chaofeng Chen, Jingwen Hou, Annan Wang, Wenxiu Sun, Qiong Yan, and Weisi Lin. Exploring Video Quality Assessment on User Generated Contents from Aesthetic and Technical Perspectives, 2023. [arXiv:2211.04894 \[cs, eess\]](https://arxiv.org/abs/2211.04894). 3
- 696
- 697
- 698 [46] Xinhao Xiang, Xiao Liu, Zizhong Li, Zhuosheng Liu, and Jiawei Zhang. Aigve-tool: Ai-generated video evaluation toolkit with multifaceted benchmark. *arXiv preprint arXiv:2503.14064*, 2025. 1, 3
- 699
- 700
- 701
- 702 [47] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. 3
- 703
- 704
- 705 [48] Boqiang Zhang, Kehan Li, Zesen Cheng, Zhiqiang Hu, Yuqian Yuan, Guanzheng Chen, Sicong Leng, Yuming Jiang, Hang Zhang, Xin Li, Peng Jin, Wenqi Zhang, Fan Wang, Lidong Bing, and Deli Zhao. Videollama 3: Frontier multimodal foundation models for image and video understanding, 2025. 5
- 706
- 707
- 708 [49] Fan Zhang, Shulin Tian, Ziqi Huang, Yu Qiao, and Ziwei Liu. Evaluation agent: Efficient and promptable evaluation framework for visual generative models. *arXiv preprint arXiv:2412.09645*, 2024. 2
- 709
- 710
- 711 [50] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert, 2020. 5, 1
- 712
- 713
- 714
- 715
- 716
- 717
- 718
- 719
- 720
- 721
- 722
- 723
- [51] Xinyi Zhang, Renyu Zhang, K. Goh, and Chenshuo Sun. The value of ai-generated metadata for ugc platforms: Evidence from a large-scale field experiment. *ArXiv*, abs/2412.18337, 2024. 1
- [52] Ming Zhong, Yang Liu, Da Yin, Yuning Mao, Yizhu Jiao, Pengfei Liu, Chenguang Zhu, Heng Ji, and Jiawei Han. Towards a unified multi-dimensional evaluator for text generation, 2022. 5, 1
- [53] Xunchu Zhou, Xiaohong Liu, Yunlong Dong, Tengchuan Kou, Yixuan Gao, Zicheng Zhang, Chunyi Li, Haoning Wu, and Guangtao Zhai. Light-vqa+: A video quality assessment model for exposure correction with vision-language guidance, 2024. 5
- 724
- 725
- 726
- 727
- 728
- 729
- 730
- 731
- 732
- 733
- 734
- 735
- 736