

# CONTENT-RICH AIGC VIDEO QUALITY ASSESSMENT VIA INTRICATE TEXT ALIGNMENT AND MOTION-AWARE CONSISTENCY

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The advent of next-generation video generation models like *Sora* poses challenges for AI-generated content (AIGC) video quality assessment (VQA). These models substantially mitigate flickering artifacts prevalent in prior models, enable longer and complex text prompts, and generate longer videos with intricate, diverse motion patterns. Conventional VQA methods designed for simple text and basic motion patterns struggle to evaluate these content-rich videos. To this end, we propose **CRAVE** (Content-Rich AIGC Video Evaluator), specifically for the evaluation of Sora-era AIGC videos. CRAVE proposes the multi-granularity text-temporal fusion that aligns long-form complex textual semantics with video dynamics. Additionally, CRAVE leverages the hybrid motion-fidelity modeling to assess temporal artifacts. Furthermore, given the straightforward prompts and content in current AIGC VQA datasets, we introduce **CRAVE-DB**, an *ITU-Compliant* benchmark featuring content-rich videos from next-generation models paired with elaborate prompts. Extensive experiments have shown that the proposed CRAVE achieves excellent results on multiple AIGC VQA benchmarks, demonstrating a high degree of alignment with human perception. All data and code will be publicly available to foster future research.

## 1 INTRODUCTION

Recently, text-driven video generation Brooks et al. (2024); Hunyuan (2024) has seen significant growth. However, evaluating these text-driven AI-generated videos presents unique and escalating challenges. These challenges primarily stem from two key issues: (1) the need for precise video-text alignment, especially with complex and lengthy text prompts; (2) the occurrence of distinct distortions that are not typically found in natural videos, such as irregular motion patterns and objects.

With the advancement of new-generation models, these challenges have become even more pronounced. These new-generation models, marked by the advent of Sora Brooks et al. (2024), offer substantial improvement in visual quality, characterized by rich details and content, such as Kling Kuaishou (2024), Gen-3-alpha Runway (2024), Vidu Shengshu (2024), etc. Compared with prior AIGC videos, they support **much longer and more intricate text (often over 200 characters), along with more complex motion patterns with longer duration (often over 5 seconds with the fps of 24)**. As illustrated in Figure 1, these rich contents impose greater demands on the evaluator’s ability to understand video dynamics and the alignment with complex textual semantics.

To address this, we introduce Content-Rich AIGC Video Evaluator (CRAVE) to assess the quality of these next-generation text-driven videos. CRAVE evaluates videos from three perspectives: It firstly considers the traditional visual harmony, like the previous Video Quality Assessment (VQA) method Wu et al. (2023a), which measures the aesthetics and distortions. Furthermore, CRAVE leverages a multi-granularity text-temporal fusion module to align the intricate texts with video dynamics. Additionally, CRAVE incorporates the hybrid motion-fidelity modeling that exploits hierarchical motion information to assess the temporal quality of the next-generation AIGC videos.

Besides, the gap between the naturalness and complexity of the latest AIGC videos and previous videos has become markedly apparent. To better assess current AIGC videos, we introduce CRAVE-DB, a content-rich AIGC VQA benchmark consisting of elaborate text-driven videos generated by

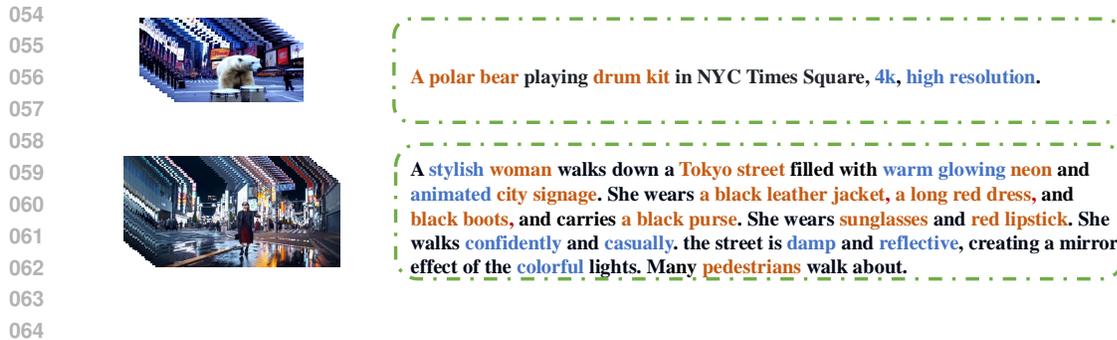


Figure 1: Comparison of concurrent and previous AIGC videos. Videos are generated by Lavie Wang et al. (2023d) (1st row) and Sora Brooks et al. (2024) (2nd row), respectively. Nouns that should be present in the video are highlighted in orange, while adjectives with more details are highlighted in blue. The new-generation AIGC videos contain richer content.

advanced models such as Kling Kuaishou (2024), Qingying Zhipu (2024), Vidu Shengshu (2024), and Sora Brooks et al. (2024). Here, by "elaborate text," we mean prompts that include complete descriptions of the subject, actions, and environment, with at least 5 detailed descriptions for any one aspect and a total character count exceeding 200. These videos have largely eliminated issues prevalent in previous generations, such as flickering, weak motion, and short content. They encompass diverse scenes, subjects, actions, and rich details, with a duration of over 5 seconds and a frame rate of 24 fps. Extensive experiments show that CRAVE has achieved leading human-aligned video quality assessment results across multiple metrics on T2V-DB Kou et al. (2024b), currently the largest AIGC VQA dataset, and the proposed CRAVE-DB.

To summarize, our main contributions are as follows: (1) We introduce CRAVE, the effective evaluator for content-rich videos derived from the new-generation video models, which assesses AIGC videos from the temporal and video-text consistency via effective motion-aware video dynamics understanding and a multi-granularity text-temporal fusion module. (2) Given the gap between new-generation AIGC videos and previous ones, we introduce CRAVE-DB, a benchmark containing AIGC VQA samples produced by advanced models like Kling, etc., to facilitate the evaluation of contemporary content-rich AIGC videos. (3) Extensive experiments demonstrate the proposed CRAVE achieves excellent results on multiple AIGC VQA benchmarks with varying sources of videos and prompt lengths, showcasing a strong understanding of the quality of AIGC videos.

## 2 RELATED WORK

### 2.1 MEASUREMENT FOR TEXT-TO-VIDEO MODELS

Currently, common methods of evaluating text-driven generated videos include some objective metrics Radford et al. (2021); Unterthiner et al. (2018); Salimans et al. (2016) and human-aligned methods Kirstain et al. (2023); Qu et al. (2024); Kou et al. (2024b). Objective metrics such as CLIP-score Radford et al. (2021) measure the mean cosine similarity between the text and each frame. IS Salimans et al. (2016) utilizes the inception feature to measure the overall quality of image and video frames. However, these objective metrics do not align with human subjective perception and often evaluate videos from a single dimension. Some measurements for natural videos provide human-aligned overall evaluations Wu et al. (2023a; 2022); Kou et al. (2023). DOVER Wu et al. (2023a) assesses quality in terms of aesthetics and technicality. FastVQA Wu et al. (2022) utilizes grid mini-patch sampling to assess videos efficiently while maintaining accuracy. Q-Align Wu et al. (2023b) transforms the VQA task into the generation of discrete quality level words via the Multimodal Large Language Model. StableVQA Chai et al. (2023) measures video stability by separately obtaining the raw optical flow, semantic, and blur features. These are suitable for natural video quality assessment but do not consider the text-video alignment, which is key to the evaluation of text-driven videos. To address this, EvalCrafter Liu et al. (2024) establishes a series of indicators including CLIP score, SD score, and natural video assessment methods. T2V-QA Kou et al. (2024b) incorporates a transformer-based encoder and a Large Language Model to assess text-driven AIGC

108 videos. TriVQA Qu et al. (2024) explores the video-text consistency through cross-attention pooling  
 109 and the recaption of Video-LLaVA. However, there are still relatively few VQA methods specifically  
 110 for AIGC videos. With the growth of new-generation videos, the requirements for understanding  
 111 video dynamics and text consistency are becoming increasingly demanding, posing greater challenges.

## 112 2.2 TEXT-TO-VIDEO GENERATION METHOD

113 Recently, lots of video generation models based on Rombach et al. (2022); Ho et al. (2020) have  
 114 emerged Singer et al. (2023); Wang et al. (2023c;a); Blattmann et al. (2023); Chen et al. (2023a);  
 115 Zheng et al. (2024); Lab & etc. (2024). They represent a significant breakthrough in video generation.  
 116 However, videos produced by prior methods still tend to suffer from issues such as low resolution,  
 117 short duration, flickering, and distortion. With the advent of Sora Brooks et al. (2024), the new-  
 118 generation models Hunyuan (2024); LumaLabs (2024); MiniMax (2024); Tongyi (2024); Labs (2024);  
 119 Yang et al. (2024) have made notable progress. Particularly recently, methods like Kling Kuaishou  
 120 (2024), Gen-3-alpha Runway (2024), and Qingying Zhipu (2024) have achieved impressive video  
 121 generation results and have been made available for community testing. These videos generally  
 122 alleviate the foundational problems seen in previous methods, with a duration of more than 5 seconds  
 123 and frame rates above 24 fps. Meanwhile, the content in these videos includes a lot of details, and they  
 124 support the control via longer text inputs. Under the wave of new-generation video generation models,  
 125 effectively assessing more complex spatiotemporal relationships within the videos and exploring their  
 126 consistency with longer texts is a topic worthy of further study.

## 127 2.3 TEXT-TO-VIDEO VQA DATASET

128 Currently, there are still relatively few  
 129 text-to-video QA datasets suitable for  
 130 evaluating current AIGC videos. Eval-  
 131 Crafter Liu et al. (2024) collects 700  
 132 prompts and uses 5 models to gener-  
 133 ate 2500 videos in total. FETV Liu  
 134 et al. (2023) utilizes 619 prompts to  
 135 generate 2,476 videos by 4 T2V mod-  
 136 els. Chivileva Chivileva et al. (2023)

137 derives 1,005 videos generated from 5 T2V models. VBench Huang et al. (2024a) uses nearly  
 138 1,700 prompts and 4 T2V models to generate 6984 videos. T2VQA-DB Kou et al. (2024a) contains  
 139 10,000 videos generated by 1000 prompts. These datasets mainly meet two challenges: (1) According  
 140 to the ITU-standard Series (2012), the number of human annotators should exceed 15 to keep the  
 141 assessment error within a controllable range. Among these, only T2VQA-DB Kou et al. (2024a) and  
 142 Chivileva Chivileva et al. (2023) meet the standard with 27 and 24 annotators. (2) The gap between  
 143 the prior and the concurrent AIGC videos. Prior videos often involve only easy movements and  
 144 commonly have basic issues such as flickering, which are relatively rarely seen in the new-generation  
 145 video models. In this work, to address the issue that prior VQA datasets do not cover concurrent AIGC  
 146 videos, we introduce CRAVE-DB, which focuses on next-generation AIGC videos with subjective  
 147 scores from **29** annotators, to provide a robust assessment of concurrent AIGC videos.

## 148 3 CONTENT-RICH AIGC VQA BENCHMARK

149 With the advancements in text-driven video generation, there exists a significant gap between the  
 150 concurrent and previous models in terms of visual quality, content complexity, and cross-modal  
 151 understanding, as shown in Figure 1. These models have substantially alleviated basic issues such  
 152 as flickering prevalent in earlier models, and have removed the prior length limitation of 77 tokens  
 153 for input text. The challenges now shift towards content distortion in more complex spatiotemporal  
 154 scenarios and semantic alignment with more intricate texts. However, current AIGC VQA datasets  
 155 are still based on the previous generation of general models, creating a significant gap compared to the  
 156 concurrent content-rich models. To this end, we introduce CRAVE-DB, a new AIGC VQA benchmark  
 157 featuring intricate text prompts, content-rich videos generated by SOTA generation models, and the  
 158 corresponding human scores. CRAVE-DB incorporates 410 intricate prompts, each containing dense  
 159 text information, as shown in Table 1. Each video has a duration of over 5 seconds with a fps of 24.

160 Table 1: Comparison of prompt density and annotations.

| Dataset                          | # Words      | # Chars.      | # Ann.    |
|----------------------------------|--------------|---------------|-----------|
| FETV Liu et al. (2023)           | 10.94        | 59.60         | 3         |
| EvalCrafter Liu et al. (2024)    | 12.33        | 69.77         | -         |
| T2V-CompBench Sun et al. (2024a) | 10.42        | 56.42         | 3         |
| VBench Huang et al. (2024b)      | 7.64         | 41.95         | -         |
| VideoGenEval Zeng et al. (2024)  | 33.00        | 202.02        | 0         |
| T2VQA-DB Kou et al. (2024b)      | 12.32        | 76.22         | 27        |
| <b>CRAVE-DB (Ours)</b>           | <b>68.38</b> | <b>411.30</b> | <b>29</b> |

In the subjective study, to ensure the data quality, we engage 29 humans to rate each video, nearly twice the minimum number of subjects required by the ITU standard Series (2012), to minimize the impact of variance in the ratings. From these, we screen out 35,612 Ratings ratings with 1,228 high-quality annotated videos to further enhance the quality of the dataset. We will subsequently introduce the prompt collection, video generation, and subjective study in detail.

### 3.1 PROMPT COLLECTION

In the past AIGC VQA datasets composed of prior-generation models, most supported prompt length is limited by CLIP Radford et al. (2021). In this case, these prompts tend to be brief, making it challenging to incorporate complex descriptions and scene compositions. For instance, we present the prompt density (average word and character count per prompt) of different datasets, as shown in Table 1. We could learn that most prompts in previous datasets contain merely a dozen words. This inherent limitation poses significant challenges for models when evaluating more sophisticated semantic alignment.

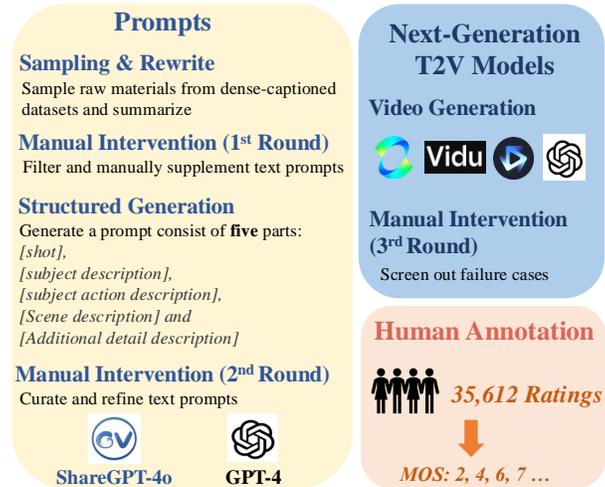


Figure 2: The collection pipeline of CRAVE-DB.

To address this, we propose the prompt generation pipeline in Figure 2. To ensure the prompts are detailed and semantically rich, we focus on the dense-captioned dataset, ShareGPT-4o Chen et al. (2023b), which leverages the advanced multimodal capabilities of GPT-4o to describe videos in detail. It contains rich annotations that even require summarization to be clear prompts. We randomly sampled 300 captions and summarized them using GPT-4 Achiam et al. (2023), retaining only key details. We then conducted the 1st intervention to filter out failed, redundant, or illogical generations.

Given that ShareGPT-4o primarily focuses on daily life scenarios, we manually crafted 200 more prompts to broaden the coverage of actions, subjects, and scenes. Prompts contain 4 categories: landscape, object, animal, and human. The "landscape" contains common scenes (e.g., grasslands, streets), rare environments (e.g., volcanoes, auroras), and renowned landmarks. The "animal" includes various mammals, reptiles, birds, fish, and amphibians. The "object" covers common real-world items, while the "human" features people across ages, genders, occupations, and clothing.

Subsequently, we employed GPT-4 to structure raw prompts using a template format: "[shot language] + [subject description] + [subject action description] + [scene description] + [additional detail description]". The "shot language" incorporates various cinematographic techniques including tilt shots, flat shots, progressive shots, surround shots, close-ups, and panoramic views. The scene descriptions encompass natural landscapes under diverse weather and lighting conditions. Following this, we initiated a second round of manual intervention to screen and refine all prompts, ultimately finalizing a curated set of 410 high-quality prompts. The overall word cloud is shown in Figure 3.

### 3.2 VIDEO GENERATION

Since the advent of Sora Brooks et al. (2024), text-driven video generation methods have achieved significant advancements in visual quality, text understanding, and the diversity and complexity of generated content. Given the substantial gap between current AIGC videos and prior ones, constructing datasets using the next-generation video models is essential. In this work, we employ Sora and other subsequent state-of-the-art models: Kling Kuaishou (2024), Vidu Shengshu (2024), Qingying Zhipu (2024) to build samples. As Sora had not been publicly released by the time of our subjective evaluation, we curated 14 content-rich prompts and their corresponding outputs from Sora's publicly showcased videos. All videos exceed 5 seconds in duration with a frame rate of 24 fps, and resolutions ranging from  $384 \times 688$  to  $960 \times 1440$ , depending on the generation model. Considering the limited availability of publicly accessible AIGC VQA datasets that meet ITU



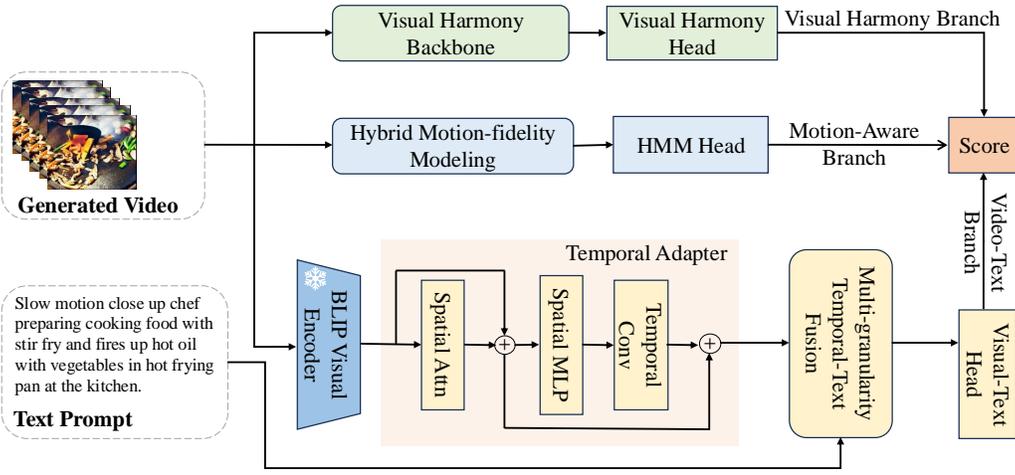


Figure 4: Network overview of the proposed CRAVE.

videos based on two dimensions: aesthetic score and technical distortion. In this work, we use the pre-trained DOVER method as the Visual Harmony Backbone, followed by a linear head to obtain the output, which could be formulated as:

$$F_{aes} = \Phi_{aes}(V), \quad (2)$$

$$F_{tech} = \Phi_{tech}(V), \quad (3)$$

$$O_{vh} = \omega_{vh}(F_{aes} \oplus F_{tech}), \quad (4)$$

where  $V$  is the input video,  $\Phi_{aes}$ ,  $\Phi_{tech}$  represent the aesthetic encoder and distortion encoder in DOVER, respectively,  $\oplus$  denotes concatenation along the dimension,  $\omega_{vh}$  represents the linear head for this branch, and  $O_{vh}$  refers to the corresponding output.

### 4.3 MULTI-GRANULARITY TEXT-TEMPORAL FUSION

Multi-granularity Text-Temporal (MTT) module firstly leverages high-quality priors from multi-modal understanding approaches like BLIP Li et al. (2022b), successfully extending it in the temporal dimension via effective temporal adapters. Then the visual information aggregated in the temporal adapter interacts with the text embedding via cross-attention. To flexibly fuse effective information from the text, we additionally perform multi-granularity aggregation on the text. Namely, in addition to the whole input, we break down the text into phrases and words of varying granularity containing different levels of semantic information via SpaCy Honnibal et al. (2020). After that, the integrated text embeddings of varying granularity are measured with the output of the visual branch, as shown in Figure 5. The entire process can be formulated as:

$$F_v = \Phi_t(\Phi_s(V)), \quad (5)$$

$$F_e = \Phi_e([enc; F_v], e), \quad (6)$$

$$F_{enc} = F_e[0, \dots], \quad (7)$$

$$F_{word} = F_e[1 :, \dots], \quad (8)$$

where  $V$ ,  $e$ ,  $F_v$  are the input video, text prompt, and the derived visual feature, respectively.  $\Phi_s$ ,  $\Phi_t$ ,  $\Phi_e$  denote the spatial encoder, temporal adapter, and text encoder, respectively.  $F_{enc}$  and  $F_{word}$  refer to the  $[enc]$  token embedding and the word token embedding, representing different levels of semantics of the textual prompt. To get the phrase-level encoding for prompts, we utilize the successful word-to-phrase module in Zhu et al. (2023). The module converts the mapping between word and phrase to the mapping between word embedding and phrase embedding, based on the position mapping between word and the corresponding embedding, which could be formulated as,

$$\{p_1, p_2 \dots, p_m\} = \Phi_w(\{w_1, w_2 \dots, w_n\}), \quad (9)$$

$$\{F_{p_1}, F_{p_2} \dots, F_{p_m}\} = \Phi_{F_w}(\{F_{w_1}, F_{w_2} \dots, F_{w_n}\}), \quad (10)$$

where  $\Phi_w$  and  $\Phi_{F_w}$  refer to the mapping between word and phrase and the mapping between word embedding and phrase embedding.  $p_i, w_i$  refers to the  $i$ -th useful phrase and word in prompt, while  $F_{p_i}, F_{w_i}$  refers to the  $i$ -th phrase embedding and word embedding, respectively. After that, we calculate the cosine distance between all text features and video features, respectively. Finally, we sum them up to get the final score.

$$O_{align} = \sum_l \cos(F_l, F_v), \tag{11}$$

where  $l$  denotes levels of granularity such as the whole paragraph, phrase, and word level.

#### 4.4 HYBRID MOTION-FIDELITY MODELING

Compared with natural videos, AIGC videos usually contain unique distortions such as irregular motions that violate the physical laws. Despite improvements in recent video generation models, low-fidelity motion remains a persistent challenge. Here, motions that defy logic, deformed motions, and motions with abnormal amplitudes are collectively referred to as "low-quality" motions. To better assess motion distortions in current AIGC videos, we propose Hybrid Motion-fidelity Modeling (HMM), which hierarchically captures motion features at different granularities. Specifically, given the success of optical flow in anomaly detection Caldelli et al. (2021); Agarwal et al. (2020), we leverage dense motion information derived from flow to capture low-level motion patterns, combined with global abstract motion information from action recognition Kay et al. (2017); Goyal et al. (2017). Section 5 and Appendix E later demonstrate the effectiveness of combining these two aspects. In practice, flow features are extracted using the pre-trained StreamFlow Sun et al. (2024c), while high-level abstract motion priors are obtained from the pre-trained Uniformer Li et al. (2023). Different branches are then fused via a feed-forward network. A linear head is then employed to process the fused results and perform regression to generate the final output.

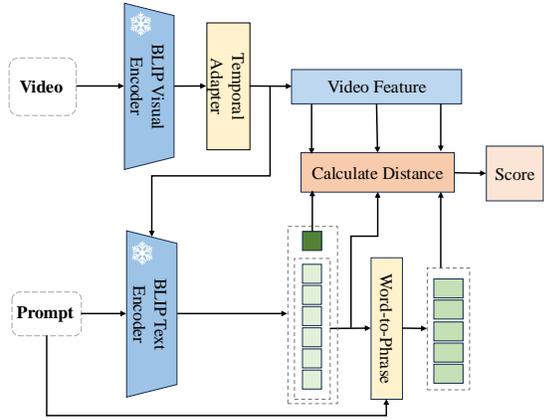


Figure 5: Details of the proposed MTT module.

#### 4.5 SUPERVISION

Based on the previous training objective in Wu et al. (2023a; 2022); Sun et al. (2022), the mix of rank loss Gao et al. (2019) and PLCC (Pearson Linear Correlation Coefficient) loss is adopted to train the overall network. Note that the BLIP visual and text encoder remains frozen throughout the training process. The whole training objective could be formulated as,

$$\mathcal{L} = \mathcal{L}_{plcc} + \gamma \cdot \mathcal{L}_{rank}, \tag{12}$$

where  $\gamma$  is the coefficient set to 0.3 in the experiments.

### 5 EXPERIMENTS

#### 5.1 EXPERIMENTAL SETUP

We use T2VQA-DB Kou et al. (2024b), GAIA Chen et al. (2024), VideoGenEval Zeng et al. (2024), and the proposed CRAVE-DB for evaluation. T2VQA-DB is the largest AIGC VQA dataset for text-driven video generation. It contains many AIGC videos from classic methods, which provides a good complement to the evaluation. We further perform cross-dataset tests on GAIA and VideoGenEval. As discussed in Section 3.2, since VideoGenEval does not have annotations, we use the Hunyuan and Wan subset of VideoGenEval as the test set and conduct subjective annotations following the ITU standard. During training and evaluation, we follow the same settings in DOVER Wu et al. (2023a).

#### 5.2 QUANTITATIVE RESULTS

As shown in Table 2 and 3, we can learn that the leading performance of CRAVE on both the proposed content-rich dataset and the T2VQA-DB that includes prior AIGC videos. On CRAVE-DB, CRAVE demonstrates a particularly significant lead, highlighting its effectiveness in evaluating next-generation AIGC videos. On T2VQA-DB, CRAVE also outperforms previous models, even surpassing LLM-based models such as Q-Align and T2VQA, which further demonstrates the effectiveness of its design. "Ft." denotes methods that need further fine-tuning on the target dataset. In Table 3, "Bg.", "Sub.", "Consis", "Aes.", "Sm." denotes the background, subject, consistency, aesthetic, and smoothness, respectively. It can be seen that 0-shot methods tend to have lower results, which is also observed in previous works Kou et al. (2024b); Sun et al. (2024b). It could be due to the lack of alignment with human perception or the consideration of the dynamic distortion in AIGC videos. We further include comparisons in GPU memory and latency in Appendix D, which demonstrate the efficiency of CRAVE.

### 5.3 QUALITATIVE RESULTS

We first visualize the difference between the predicted and the ground-truth MOS, as illustrated in Appendix F. The curves are obtained from a fourth-order polynomial nonlinear fitting. We then further present cases of CRAVE-DB and limitations of CRAVE in Appendix F and G, which include the scores of various samples by different quality assessment models and failure cases of CRAVE, respectively. It can be seen that models trained on CRAVE-DB demonstrate more accurate evaluations for those cases compared to training on T2VQA-DB, indicating that CRAVE-DB serves as an effective complement to existing datasets. More video cases could be found in the **supplemental materials**.

### 5.4 CROSS-DATASET VALIDATION

In this section, all models are trained on T2VQA-DB given its large capacity. They are then performed 0-shot testing on CRAVE-DB, GAIA, and VideoGenEval, as shown in Table 4. We choose GAIA as it is one of a limited number of public ITU-Compliant AIGC VQA datasets (Annotators no less than 15). We chose VideoGenEval based on the discussion in Section 3.2. Since VideoGenEval does not have labels, we conduct annotations based on settings in Section 3.3. Considering the annotation cost and the distinctiveness from other datasets, we selected videos from the latest video generation models Wan and Hunyuan subsets, which include 465 samples. It can be seen that CRAVE demonstrates excellent cross-dataset generalization capability across different test sets.

### 5.5 ABLATION STUDY

We further ablate each component of CRAVE, as shown in Table 5. Underlined settings are used in our final model. As CRAVE-DB naturally contains intricate texts, rich motion information, and other

Table 2: Quantitative comparison on T2VQA-DB.

| Type   | Models                        | SRCC          | PLCC          | KRCC          |
|--------|-------------------------------|---------------|---------------|---------------|
| 0-shot | CLIPSim Radford et al. (2021) | 0.1047        | 0.1277        | 0.0702        |
|        | BLIP Li et al. (2022b)        | 0.1659        | 0.1860        | 0.1112        |
|        | ImageReward Xu et al. (2023)  | 0.1875        | 0.2121        | 0.1266        |
|        | ViCLIP Wang et al. (2023e)    | 0.1162        | 0.1449        | 0.0781        |
|        | UMTScore Liu et al. (2023)    | 0.0676        | 0.0721        | 0.0453        |
| Ft.    | SimpleVQA Sun et al. (2022)   | 0.6275        | 0.6338        | 0.4466        |
|        | BVQA Li et al. (2022a)        | 0.7390        | 0.7486        | 0.5487        |
|        | FastVQA Wu et al. (2022)      | 0.7173        | 0.7295        | 0.5303        |
|        | DOVER Wu et al. (2023a)       | 0.7609        | 0.7693        | 0.5704        |
|        | Q-Align Wu et al. (2023b)     | 0.7601        | 0.7768        | 0.5860        |
|        | T2VQA Kou et al. (2024b)      | 0.7965        | 0.8066        | 0.6058        |
| Ours   | CRAVE                         | <b>0.8122</b> | <b>0.8214</b> | <b>0.6338</b> |

Table 3: Quantitative comparison on CRAVE-DB.

| Type                        | Models                            | SRCC                        | PLCC          | KRCC          |
|-----------------------------|-----------------------------------|-----------------------------|---------------|---------------|
| 0-shot                      | UMTScore Liu et al. (2023)        | 0.0134                      | 0.0355        | 0.0118        |
|                             | Flicker Huang et al. (2024b)      | 0.0761                      | 0.0687        | 0.0493        |
|                             | Bg. Consis. Huang et al. (2024b)  | 0.1170                      | 0.0829        | 0.0773        |
|                             | Sub. Consis. Huang et al. (2024b) | 0.1418                      | 0.0956        | 0.0936        |
|                             | ViCLIP Wang et al. (2023e)        | 0.1290                      | 0.1207        | 0.0857        |
|                             | LongCLIP Zhang et al. (2024)      | 0.2757                      | 0.3033        | 0.1884        |
|                             | Aes. Score Huang et al. (2024b)   | 0.3557                      | 0.3457        | 0.2373        |
|                             | HPSv2 Wu et al. (2023c)           | 0.0562                      | 0.0501        | 0.0391        |
|                             | PickScore Kirstain et al. (2023)  | 0.0557                      | 0.0492        | 0.0361        |
|                             | CHScore Yuan et al. (2024)        | 0.2202                      | 0.1628        | 0.1470        |
|                             | ImageReward Xu et al. (2023)      | 0.2216                      | 0.2125        | 0.1470        |
|                             | Motion Sm. Huang et al. (2024b)   | 0.2331                      | 0.2630        | 0.1549        |
|                             | Ft.                               | SimpleVQA Sun et al. (2022) | 0.6230        | 0.6180        |
| StableVQA Kou et al. (2023) |                                   | 0.6396                      | 0.6415        | 0.4552        |
| Q-Align Wu et al. (2023b)   |                                   | 0.6481                      | 0.6492        | 0.4712        |
| FastVQA Wu et al. (2022)    |                                   | 0.7155                      | 0.7062        | 0.5243        |
| DOVER Wu et al. (2023a)     |                                   | 0.7095                      | 0.7192        | 0.5224        |
| TriVQA Qu et al. (2024)     |                                   | 0.7183                      | 0.7313        | 0.5293        |
| T2VQA Kou et al. (2024b)    | 0.7266                            | 0.7098                      | 0.5369        |               |
| Ours                        | CRAVE                             | <b>0.7587</b>               | <b>0.7581</b> | <b>0.5660</b> |

Table 4: Cross-dataset validation on CRAVE-DB, GAIA, and VideoGenEval.

| Method                    | CRAVE-DB      |               | GAIA          |               | VGenEval-Hunyuan |               | VGenEval-Wan  |               |
|---------------------------|---------------|---------------|---------------|---------------|------------------|---------------|---------------|---------------|
|                           | SRCC          | PLCC          | SRCC          | PLCC          | SRCC             | PLCC          | SRCC          | PLCC          |
| DOVER Wu et al. (2023a)   | 0.3480        | 0.3719        | 0.1187        | 0.1141        | 0.4363           | 0.3663        | 0.3284        | 0.3259        |
| Q-Align Wu et al. (2023b) | 0.3757        | 0.3943        | 0.2105        | 0.1747        | 0.4331           | 0.4335        | 0.2943        | 0.3013        |
| T2VQA Kou et al. (2024b)  | 0.4064        | 0.4166        | 0.2198        | 0.2166        | 0.4340           | 0.4498        | 0.3570        | 0.3785        |
| CRAVE (Ours)              | <b>0.4213</b> | <b>0.4326</b> | <b>0.2263</b> | <b>0.2258</b> | <b>0.4595</b>    | <b>0.4694</b> | <b>0.3853</b> | <b>0.3805</b> |

Table 5: Quantitative results of ablation study.

| Experiment                       | Method                | CRAVE-DB      |               |               | T2VQA-DB      |               |               |
|----------------------------------|-----------------------|---------------|---------------|---------------|---------------|---------------|---------------|
|                                  |                       | SRCC          | PLCC          | KRCC          | SRCC          | PLCC          | KRCC          |
| Temporal-Text Fusion             | None                  | 0.6512        | 0.6601        | 0.4643        | 0.7701        | 0.7804        | 0.5885        |
|                                  | ST-Graph.             | 0.6966        | 0.6968        | 0.5042        | 0.7990        | 0.8084        | 0.6196        |
|                                  | Temp. Attn.           | 0.7159        | 0.7175        | 0.5226        | 0.7989        | 0.8089        | 0.6199        |
|                                  | <u>Pseudo 3D Conv</u> | <b>0.7310</b> | <b>0.7241</b> | <b>0.5335</b> | <b>0.8077</b> | <b>0.8167</b> | <b>0.6289</b> |
| Multi-Granularity Text Injection | None                  | 0.7477        | 0.7458        | 0.5526        | 0.8077        | 0.8167        | 0.6289        |
|                                  | Global                | 0.7520        | 0.7485        | 0.5566        | 0.8115        | 0.8207        | 0.6329        |
|                                  | Local                 | 0.7524        | 0.7498        | 0.5578        | 0.8121        | 0.8212        | 0.6337        |
|                                  | Individual            | 0.7491        | 0.7460        | 0.5534        | 0.8111        | 0.8200        | 0.6326        |
|                                  | <u>Combined</u>       | <b>0.7587</b> | <b>0.7581</b> | <b>0.5660</b> | <b>0.8122</b> | <b>0.8214</b> | <b>0.6338</b> |
| Motion-Aware Modeling            | None                  | 0.7507        | 0.7499        | 0.5561        | 0.8103        | 0.8199        | 0.6317        |
|                                  | High-Level            | 0.7527        | 0.7529        | 0.5605        | 0.8111        | 0.8204        | 0.6326        |
|                                  | Low-Level             | 0.7529        | 0.7513        | 0.5576        | 0.8118        | 0.8211        | 0.6330        |
|                                  | <u>Hybrid</u>         | <b>0.7587</b> | <b>0.7581</b> | <b>0.5660</b> | <b>0.8122</b> | <b>0.8214</b> | <b>0.6338</b> |
| Flow Frames                      | 4                     | 0.7562        | 0.7552        | 0.5598        | 0.8110        | 0.8206        | 0.6323        |
|                                  | 8                     | 0.7566        | 0.7564        | 0.5612        | 0.8113        | 0.8208        | 0.6326        |
|                                  | <u>16</u>             | <b>0.7587</b> | <b>0.7581</b> | <b>0.5660</b> | <b>0.8122</b> | <b>0.8214</b> | <b>0.6338</b> |

such content, we could learn that the improvements on it are generally more significant. We first explore ways to align the text with temporal visual features. ST-Graph denotes the spatio-temporal graph modeling, which flattens the temporal dimension into the spatial. Temp. Attn. is the attention operator along the additional temporal dimension. Pseudo 3D Conv Singer et al. (2023) stacks additional convolutions in the temporal dimension. We can see that the effectiveness has significantly improved with temporal modeling, and that the Pseudo 3D Conv widely used in generation tasks also excels in long-text spatiotemporal modeling. We then investigate the granularity of MTT and discover that integrating all granularity levels yields optimal performance. A more detailed analysis for each branch is in Appendix E. Additionally, we examine the impact of motion-aware temporal modeling. Our experiments demonstrate that dense data from optical flow enhances overall performance, and incorporating sparse abstract spatiotemporal information provides a significant performance boost. We further explored the impact of flow frames. We observed that using more optical flow frames tends to improve accuracy. Given the trade-off between accuracy and efficiency, we ultimately chose 16 frames for the flow calculation. Besides, we ablate the effect of various backbones in visual harmony, visual-text alignment, and motion-aware modeling. Please refer to Appendix E for details.

## 6 CONCLUSION

Given the gap between concurrent text-driven video generation models and the existing AIGC VQA dataset, we introduce CRAVE, an effective VQA method, and CRAVE-DB, a new benchmark for the next-generation AIGC videos. Based on the effective multi-dimensional design, CRAVE achieves excellent human-aligned results across multiple metrics and datasets. CRAVE-DB better aligns with contemporary AIGC-generated videos, serving as an effective complement to existing datasets.

## REFERENCES

- 486  
487  
488 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
489 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.  
490 *arXiv preprint arXiv:2303.08774*, 2023.
- 491 Shruti Agarwal, Hany Farid, Tarek El-Gaaly, and Ser-Nam Lim. Detecting deep-fake videos from  
492 appearance and behavior. In *2020 IEEE international workshop on information forensics and*  
493 *security (WIFS)*, pp. 1–6. IEEE, 2020.
- 494 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik  
495 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling  
496 latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- 497  
498 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe  
499 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video  
500 generation models as world simulators. Technical report, OpenAI, 2024.
- 501 Roberto Caldelli, Leonardo Galteri, Irene Amerini, and Alberto Del Bimbo. Optical flow based cnn  
502 for detection of unlearned deepfake manipulations. *Pattern Recognition Letters*, 146:31–37, 2021.
- 503  
504 Wenhao Chai, Xun Guo, Gaoang Wang, and Yan Lu. Stablevideo: Text-driven consistency-aware  
505 diffusion video editing. In *Proceedings of the IEEE/CVF International Conference on Computer*  
506 *Vision*, pp. 23040–23050, 2023.
- 507 Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo  
508 Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, et al. Videocrafter1: Open diffusion models for  
509 high-quality video generation. *arXiv preprint arXiv:2310.19512*, 2023a.
- 510  
511 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong  
512 Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. Internvl:  
513 Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint*  
514 *arXiv:2312.14238*, 2023b.
- 515 Zijian Chen, Wei Sun, Yuan Tian, Jun Jia, Zicheng Zhang, Wang Jiarui, Ru Huang, Xiongkuo Min,  
516 Guangtao Zhai, and Wenjun Zhang. Gaia: Rethinking action quality assessment for ai-generated  
517 videos. *Advances in Neural Information Processing Systems*, 37:40111–40144, 2024.
- 518 Iya Chivileva, Philip Lynch, Tomas E Ward, and Alan F Smeaton. Measuring the quality of text-to-  
519 video model outputs: Metrics and dataset. *arXiv preprint arXiv:2309.08009*, 2023.
- 520  
521 Fei Gao, Dacheng Tao, Xinbo Gao, and Xuelong Li. Learning to rank for blind image quality  
522 assessment, 2019.
- 523 Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal,  
524 Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The  
525 “something something” video database for learning and evaluating visual common sense. In  
526 *Proceedings of the IEEE international conference on computer vision*, pp. 5842–5850, 2017.
- 527  
528 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
529 *neural information processing systems*, 33:6840–6851, 2020.
- 530 Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. spacy: Industrial-  
531 strength natural language processing in python. 2020. doi: 10.5281/zenodo.1212303.
- 532  
533 Ziqi Huang, Yanan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing  
534 Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin,  
535 Yu Qiao, and Ziwei Liu. Vbench: Comprehensive benchmark suite for video generative models.  
536 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024a.
- 537  
538 Ziqi Huang, Yanan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing  
539 Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video  
generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
*Recognition*, pp. 21807–21818, 2024b.

- 540 Tencent Hunyuan. Hunyuanvideo: A systematic framework for large video generative models, 2024.  
541 URL <https://arxiv.org/abs/2412.03603>.
- 542
- 543 Int.Telecommun.Union. Methodology for the subjective assessment of the quality of television  
544 pictures itu-t recommendation. *Tech. Rep.*, 2000.
- 545 Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan,  
546 Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset.  
547 *arXiv preprint arXiv:1705.06950*, 2017.
- 548
- 549 Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-  
550 a-pic: An open dataset of user preferences for text-to-image generation. *Advances in Neural  
551 Information Processing Systems*, 36:36652–36663, 2023.
- 552 Tengchuan Kou, Xiaohong Liu, Wei Sun, Jun Jia, Xiongkuo Min, Guangtao Zhai, and Ning Liu.  
553 Stablevqa: A deep no-reference quality assessment model for video stability. In *Proceedings of the  
554 31st ACM International Conference on Multimedia*, pp. 1066–1076, 2023.
- 555 Tengchuan Kou, Xiaohong Liu, Zicheng Zhang, Chunyi Li, Haoning Wu, Xiongkuo Min, Guangtao  
556 Zhai, and Ning Liu. Subjective-aligned dataset and metric for text-to-video quality assessment,  
557 2024a. URL <https://arxiv.org/abs/2403.11956>.
- 558
- 559 Tengchuan Kou, Xiaohong Liu, Zicheng Zhang, Chunyi Li, Haoning Wu, Xiongkuo Min, Guangtao  
560 Zhai, and Ning Liu. Subjective-aligned dataset and metric for text-to-video quality assessment.  
561 *arXiv preprint arXiv:2403.11956*, 2024b.
- 562 Kuaishou. Kling. <https://kling.kuaishou.com/>, 2024.
- 563
- 564 PKU-Yuan Lab and Tuzhan AI etc. Open-sora-plan, April 2024. URL [https://doi.org/10.  
565 5281/zenodo.10948109](https://doi.org/10.5281/zenodo.10948109).
- 566 Pika Labs. Pika 1.5. <https://pika.art>, 2024.
- 567
- 568 Bowen Li, Weixia Zhang, Meng Tian, Guangtao Zhai, and Xianpei Wang. Blindly assess quality of  
569 in-the-wild videos via quality-aware pre-training and motion perception. *IEEE Transactions on  
570 Circuits and Systems for Video Technology*, 32(9):5944–5958, 2022a.
- 571 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image  
572 pre-training for unified vision-language understanding and generation. In *ICML*, 2022b.
- 573
- 574 Kunchang Li, Yali Wang, Junhao Zhang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and  
575 Yu Qiao. Uniformer: Unifying convolution and self-attention for visual recognition. *IEEE  
576 Transactions on Pattern Analysis and Machine Intelligence*, 45(10):12581–12600, 2023.
- 577 Yaofang Liu, Xiaodong Cun, Xuebo Liu, Xintao Wang, Yong Zhang, Haoxin Chen, Yang Liu,  
578 Tieyong Zeng, Raymond Chan, and Ying Shan. Evalcrafter: Benchmarking and evaluating large  
579 video generation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and  
580 Pattern Recognition*, pp. 22139–22149, 2024.
- 581 Yuanxin Liu, Lei Li, Shuhuai Ren, Rundong Gao, Shicheng Li, Sishuo Chen, Xu Sun, and Lu Hou.  
582 Fetv: A benchmark for fine-grained evaluation of open-domain text-to-video generation. *arXiv  
583 preprint arXiv: 2311.01813*, 2023.
- 584
- 585 LumaLabs. Dream machine. <https://lumalabs.ai/dream-machine>, 2024.
- 586
- 587 MiniMax. Hailuo ai. <https://hailuoai.com/video>, 2024.
- 588
- 589 Bowen Qu, Xiaoyu Liang, Shangkun Sun, and Wei Gao. Exploring aigc video quality: A fo-  
590 cus on visual harmony, video-text consistency and domain distribution gap. *arXiv preprint  
591 arXiv:2404.13573*, 2024.
- 592
- 593 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
models from natural language supervision. In *International conference on machine learning*, pp.  
8748–8763. PMLR, 2021.

- 594 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
595 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference*  
596 *on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 597 Runway. Gen-3. <https://runwayml.com/blog/introducing-gen-3-alpha/>, 2024.
- 599 Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen.  
600 Improved techniques for training gans. *Advances in neural information processing systems*, 29,  
601 2016.
- 602 B Series. Methodology for the subjective assessment of the quality of television pictures. *Recommendation*  
603 *ITU-R BT*, 500(13), 2012.
- 605 Shengshu. Vidu. <https://www.vidu.studio/create>, 2024.
- 607 M Six Silberman, Bill Tomlinson, Rochelle LaPlante, Joel Ross, Lilly Irani, and Andrew Zaldivar.  
608 Responsible research with crowds: pay crowdworkers at least minimum wage. *Communications of*  
609 *the ACM*, 61(3):39–41, 2018.
- 610 Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry  
611 Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video  
612 data. In *The Eleventh International Conference on Learning Representations*, 2023.
- 613 Kaiyue Sun, Kaiyi Huang, Xian Liu, Yue Wu, Zihan Xu, Zhenguo Li, and Xihui Liu. T2v-  
614 compbench: A comprehensive benchmark for compositional text-to-video generation. *arXiv*  
615 *preprint arXiv:2407.14505*, 2024a.
- 617 Shangkun Sun, Xiaoyu Liang, Songlin Fan, Wenxu Gao, and Wei Gao. Ve-bench: Subjective-aligned  
618 benchmark suite for text-driven video editing quality assessment. In *Proceedings of the AAAI*  
619 *Conference on Artificial Intelligence*, 2024b.
- 620 Shangkun Sun, Jiaming Liu, Thomas H Li, Huaxia Li, Guoqing Liu, and Wei Gao. Streamflow:  
621 Streamlined multi-frame optical flow estimation for video sequences. In *Advances in neural*  
622 *information processing systems*, 2024c.
- 624 Wei Sun, Xionghuo Min, Wei Lu, and Guangtao Zhai. A deep learning based no-reference quality  
625 assessment model for ugc videos. In *Proceedings of the 30th ACM International Conference on*  
626 *Multimedia*, pp. 856–865, 2022.
- 627 Luming Tang, Menglin Jia, Qianqian Wang, Cheng Perng Phoo, and Bharath Hariharan. Emergent  
628 correspondence from image diffusion. *Advances in Neural Information Processing Systems*, 36:  
629 1363–1389, 2023.
- 630 Ali Tongyi. Wanxiang video. [https://tongyi.aliyun.com/wanxiang/](https://tongyi.aliyun.com/wanxiang/videoCreation)  
631 [videoCreation](https://tongyi.aliyun.com/wanxiang/videoCreation), 2024.
- 633 Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and  
634 Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *arXiv*  
635 *preprint arXiv:1812.01717*, 2018.
- 636 Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao,  
637 Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang, Jixuan  
638 Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi Zhang, Pandeng  
639 Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing Wang,  
640 Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng Zhou, Wentu  
641 Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yangyu  
642 Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu,  
643 Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan  
644 Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models. *arXiv preprint*  
645 *arXiv:2503.20314*, 2025.
- 646 Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Mod-  
647 elscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023a.

- 648 Rui Wang, Dongdong Chen, Zuxuan Wu, Yinpeng Chen, Xiyang Dai, Mengchen Liu, Lu Yuan, and  
649 Yu-Gang Jiang. Masked video distillation: Rethinking masked feature modeling for self-supervised  
650 video representation learning. In *Proceedings of the IEEE/CVF conference on computer vision  
651 and pattern recognition*, pp. 6312–6322, 2023b.
- 652 Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan  
653 He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent  
654 diffusion models. *arXiv preprint arXiv:2309.15103*, 2023c.
- 655 Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan  
656 He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent  
657 diffusion models. *arXiv preprint arXiv:2309.15103*, 2023d.
- 658 Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Jian Ma, Xinyuan Chen, Yaohui  
659 Wang, Ping Luo, Ziwei Liu, Yali Wang, Limin Wang, and Y. Qiao. Internvid: A large-scale  
660 video-text dataset for multimodal understanding and generation. *ArXiv*, abs/2307.06942, 2023e.  
661 URL <https://api.semanticscholar.org/CorpusID:259847783>.
- 662 Yi Han Wang, Lahav Lipson, and Jia Deng. Sea-raft: Simple, efficient, accurate raft for optical flow.  
663 In *European Conference on Computer Vision*, pp. 36–54. Springer, 2024.
- 664 Haoning Wu, Chaofeng Chen, Jingwen Hou, Liang Liao, Annan Wang, Wenxiu Sun, Qiong Yan, and  
665 Weisi Lin. Fast-vqa: Efficient end-to-end video quality assessment with fragment sampling. In  
666 *European conference on computer vision*, pp. 538–554. Springer, 2022.
- 667 Haoning Wu, Erli Zhang, Liang Liao, Chaofeng Chen, Jingwen Hou, Annan Wang, Wenxiu Sun,  
668 Qiong Yan, and Weisi Lin. Exploring video quality assessment on user generated contents from  
669 aesthetic and technical perspectives. In *Proceedings of the IEEE/CVF International Conference on  
670 Computer Vision*, pp. 20144–20154, 2023a.
- 671 Haoning Wu, Zicheng Zhang, Weixia Zhang, Chaofeng Chen, Chunyi Li, Liang Liao, Annan Wang,  
672 Erli Zhang, Wenxiu Sun, Qiong Yan, Xiongkuo Min, Guangtai Zhai, and Weisi Lin. Q-align:  
673 Teaching Imms for visual scoring via discrete text-defined levels. *arXiv preprint arXiv:2312.17090*,  
674 2023b.
- 675 Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li.  
676 Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image  
677 synthesis. *arXiv preprint arXiv:2306.09341*, 2023c.
- 678 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong.  
679 Imagereward: Learning and evaluating human preferences for text-to-image generation, 2023.
- 680 Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang,  
681 Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models  
682 with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.
- 683 Shenghai Yuan, Jinfa Huang, Yongqi Xu, Yaoyang Liu, Shaofeng Zhang, Yujun Shi, Rui-Jie Zhu,  
684 Xinhua Cheng, Jiebo Luo, and Li Yuan. Chronomagic-bench: A benchmark for metamorphic  
685 evaluation of text-to-time-lapse video generation. *Advances in Neural Information Processing  
686 Systems*, 37:21236–21270, 2024.
- 687 Ailing Zeng, Yuhang Yang, Weidong Chen, and Wei Liu. The dawn of video generation: Preliminary  
688 explorations with sora-like models. *arXiv preprint arXiv:2410.05227*, 2024.
- 689 Beichen Zhang, Pan Zhang, Xiaoyi Dong, Yuhang Zang, and Jiaqi Wang. Long-clip: Unlocking the  
690 long-text capability of clip. *arXiv preprint arXiv:2403.15378*, 2024.
- 691 Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou,  
692 Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all, March  
693 2024. URL <https://github.com/hpcaitech/Open-Sora>.
- 694 Zhipu. Ying. <https://chatglm.cn/video>, 2024.
- 695 Cunjuan Zhu, Qi Jia, Wei Chen, Yanming Guo, and Yu Liu. Deep learning for video-text retrieval: a  
696 review. *International Journal of Multimedia Information Retrieval*, 12(1):3, 2023.

## A ETHICS STATEMENT

This work is conducted in strict adherence to the ICLR Code of Ethics. We affirm our commitment to contributing to societal and human well-being by focusing our research on a socially beneficial task. The utilized datasets and models involves no human subjects directly, and thus carries no risks to personal privacy, health, or safety. There is no discrimination or negative societal impact in the subjective or objective experimental phases, and all findings will be accurately reported.

## B REPRODUCIBILITY STATEMENT

To ensure reproducibility, this work provides detailed descriptions of the experimental procedures. The proposed dataset, code, and model weights used will be made publicly available. In Section 3 and Appendix H, we elaborate on the dataset construction and the details of the subjective experiments, while Section 5 specifies the model parameters. The T2VQA-DB dataset utilized is publicly accessible. We believe this information will facilitate the reproduction of our work and help advance the field.

## C USE OF LLMs

Large Language Models (LLMs) were used in this work solely for typo checking and language polishing. It is important to note that LLMs were not involved in the ideation, research methodology, or experiment design of this work. All research concepts, ideas, experiments, and analyses were conducted by the authors. The authors take full responsibility for the content of the manuscript, including the text polished by LLMs. We have verified that the LLM-generated text adheres to ethical guidelines and is free from issues such as plagiarism or scientific misconduct.

## D COMPARISON ON MODEL EFFICIENCY

In this section, we compare the computational costs (inference speed, GPU memory usage, and parameter count) of different quality assessment models, as shown in Table 6. The original video input resolution is 1280x720, and the tests were conducted on NVIDIA RTX 3090 GPUs. The reported time is the average result of 5 measurements. From the table, it could be seen that CRAVE, achieves superior alignment performance compared to methods with equal or even higher parameter counts and more powerful multimodal foundation models. Furthermore, CRAVE does not fall behind in terms of inference speed, memory usage, and parameter efficiency.

Table 6: Comparison on inference speed, GPU memory usage, and parameter count.

| Method       | GPU Memory  | Parameter Counts | FPS         | SRCC          | PLCC          |
|--------------|-------------|------------------|-------------|---------------|---------------|
| Dover        | <b>3.5G</b> | <b>58.1M</b>     | <b>63.5</b> | 0.7095        | 0.7192        |
| StableVQA    | 6.1G        | 118.0 M          | 61.8        | 0.6396        | 0.6415        |
| TriVQA       | 7.7G        | 253.8M           | 56.3        | 0.7183        | 0.7103        |
| Q-Align      | 21.5G       | 5B               | 21.8        | 0.6481        | 0.6492        |
| T2VQA        | 16.2G       | 7B               | 36.4        | 0.7266        | 0.7098        |
| CRAVE (Ours) | 10.8G       | 731.7M           | 50.1        | <b>0.7587</b> | <b>0.7581</b> |

## E ADDITIONAL ABLATED RESULTS

In this section, we first further explore the impact of each branch in MTT, as shown in Table 7. We observe that varying semantic granularity brought by a single branch can influence the vision-language alignment differently. For example, as shown in the first two rows, replacing Local + Individual branches with Local + Local, the performance will drop due to the absence of individual-level granularity. If we further replace the Global branch that provides global information, with the

Individual branch, and downgrade Local information to Individual information, the performance further drops, even though these experiments have the same number of branches and parameters. This highlights the necessity of multi-granularity.

Table 7: Ablation study on MTT branches.

| Method                               | SRCC   | PLCC   |
|--------------------------------------|--------|--------|
| Global + Location + Individual       | 0.7587 | 0.7581 |
| Global + Location + Location         | 0.7554 | 0.7567 |
| Individual + Location + Location     | 0.7552 | 0.7527 |
| Individual + Individual + Individual | 0.7435 | 0.7473 |

To better illustrate how hybrid motion modeling works, we further divide the CRAVE dataset. Following EvalCrafter Liu et al. (2024), we use Flow-Score (average optical flow magnitude of videos) to measure video dynamics and split the test set into low-dynamic (Flow-Score less than 5) and high-dynamic (Flow-Score  $\geq$  5) subsets. The performance of different modeling methods is presented in Table 8. We find that CRAVE shows excellent performance even in highly dynamic regions. Meanwhile, low-level encoders like optical flow are more sensitive to higher dynamic videos and thus perform better on the high set. In comparison, high-level action recognition backbones might emphasize action semantics, resulting in relatively less improvement in higher dynamics compared to the low-level encoder. It could also be seen that a hybrid approach combining both methods effectively leverages their complementary strengths.

Table 8: Ablation study on subset of CRAVE-DB.

| Method          | Low Set       |               | High Set      |               |
|-----------------|---------------|---------------|---------------|---------------|
|                 | SRCC          | PLCC          | SRCC          | PLCC          |
| High-Level only | 0.7301        | 0.7232        | 0.7589        | 0.7650        |
| Low-Level only  | 0.7241        | 0.7197        | 0.7914        | 0.7941        |
| Hybrid          | <b>0.7322</b> | <b>0.7231</b> | <b>0.7919</b> | <b>0.8011</b> |

We then further verify the choice of different backbones in visual harmony, visual-text alignment, high-level and low-level motion modeling, as shown in Table 9. It could be seen that the quality of the pre-trained model influences the final results within a certain range, but does not cause particularly significant fluctuations in performance.

Table 9: Ablation study on various backbones.

| Experiment                 | Method                        | SRCC          | PLCC          |
|----------------------------|-------------------------------|---------------|---------------|
| Visual Harmony Backbone    | SimpleVQA Sun et al. (2022)   | 0.7519        | 0.7485        |
|                            | TriVQA Qu et al. (2024)       | 0.7547        | 0.7556        |
|                            | DOVER Wu et al. (2023a)       | <b>0.7587</b> | <b>0.7581</b> |
| Visual-Text Alignment      | CLIP Radford et al. (2021)    | 0.7460        | 0.7302        |
|                            | DIFT Tang et al. (2023)       | 0.7347        | 0.7365        |
|                            | BLIP Li et al. (2022b)        | <b>0.7497</b> | <b>0.7501</b> |
| High-Level Motion Modeling | MVD Wang et al. (2023b)       | 0.7470        | 0.7536        |
|                            | Uniformer Li et al. (2023)    | <b>0.7587</b> | <b>0.7581</b> |
| Low-level Motion Modeling  | Sea-RAFT Wang et al. (2024)   | 0.7540        | 0.7527        |
|                            | StreamFlow Sun et al. (2024c) | <b>0.7587</b> | <b>0.7581</b> |

We also ablate the effect of the initial gamma values introduced in Section 4.5, as shown in Table 10. The initial gamma value was set following DOVER Wu et al. (2023a). We tested other values to further validate our choice. It can be seen that gamma moderately affects performance.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

Table 10: Ablation study on gamma values.

| Gamma | SRCC          | PLCC          |
|-------|---------------|---------------|
| 0.1   | 0.7583        | 0.7506        |
| 0.3   | <b>0.7587</b> | <b>0.7581</b> |
| 0.5   | 0.7507        | 0.7479        |

## F ADDITIONAL QUALITATIVE RESULTS

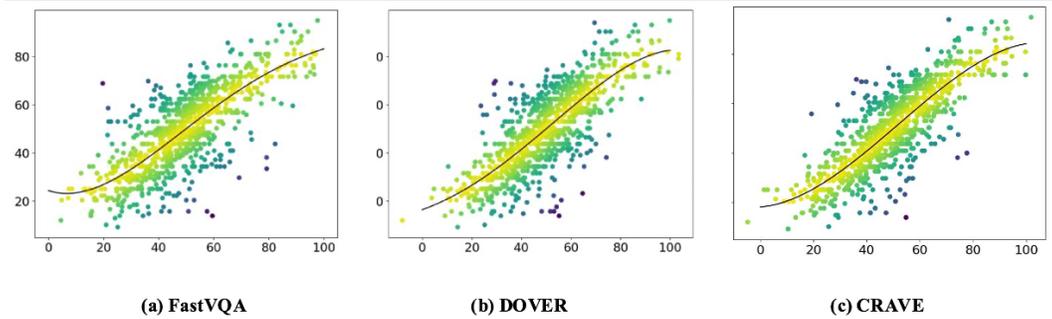


Figure 6: Scatter plots of the predicted scores and ground-truth MOSs. A brighter scatter point represents higher density.

In this section, we first visualize the difference between predicted and ground-truth MOS. As shown in Figure 6, (a), (b), and (c) represent the visualization of the differences between different models on T2VQA-DB. The more clustered the points, the smaller the differences. We can observe that the points in (a) and (b) are more scattered and farther from the central line. A fourth-order polynomial nonlinear fitting is used to draw the central line.

We further present the challenges of existing methods on the constructed dataset to show the necessity of CRAVE-DB. As shown in Figure 7, it provides detailed failure cases of multiple previous methods on the new dataset. We found that these cases often require more complex semantic alignment or occur under more dynamic conditions. Even state-of-the-art methods like T2VQA and classic approaches like Dover tend to make mistakes, suggesting that following methods could further improve consistency in complex text and motion alignment. However, after training on CRAVE-DB, **these biases are significantly reduced, which demonstrates CRAVE-DB effectively complements T2VQA-DB.**

## G LIMITATIONS AND FAILURE CASES

Overall, CRAVE also has its own limitations. First, despite CRAVE’s multiple rounds of manual screening and the use of an above-average number of annotators to ensure the diversity and reliability of the data, its applicability still faces challenges when compared to complex real-world scenarios. Second, given the current network design paradigm, the number of frames and resolution used by CRAVE during training are still limited, which may still create a gap when compared to actual videos. We believe that the evaluation of text-driven video generation tasks still has many pressing issues that need to be resolved. Although CRAVE has made some improvements compared to previous methods, there are still many problems waiting to be addressed in follow-up work. As shown in Figure 8, CRAVE still has room for improvement in evaluating long textual descriptions and local distortions. For some physical activities, the assessments made by CRAVE may still have biases in certain situations.

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917



Figure 7: Samples and corresponding scores from models trained on different datasets.



Figure 8: Failure cases of CRAVE.

## H DETAILS OF SUBJECTIVE STUDY

In this section, we give a more detailed discussion of the subjective experiments. The subjective experiments were conducted in two batches. The software interface used in the subjective experiments is illustrated in Figure 9. During the scoring, the annotators were advised to assess the video from 3 perspectives: (1) visual quality, which is commonly used in traditional VQA methods; (2) matching between the text and video; and (3) motion quality, such as whether motion consistency is maintained, whether the motion is distorted, and whether it aligns with common sense. The first batch involved annotating the CRAVE-DB dataset, with 29 participants. In this batch, all videos were divided into 10 sessions. The second batch focused on annotating the test set of VideoGenEval, involving 31 participants. All videos in this batch were divided into 5 sessions. On average, each participant spent approximately 0.91 hours completing one session. Each participant was compensated \$8.2 per session in accordance with current ethical standards Silberman et al. (2018). The experiments contained no NSFW (Not Safe For Work) or other inappropriate content in the text prompts or generated outputs. Participants came from diverse academic backgrounds and were trained using out-of-dataset cases prior to formal scoring to ensure annotation consistency.

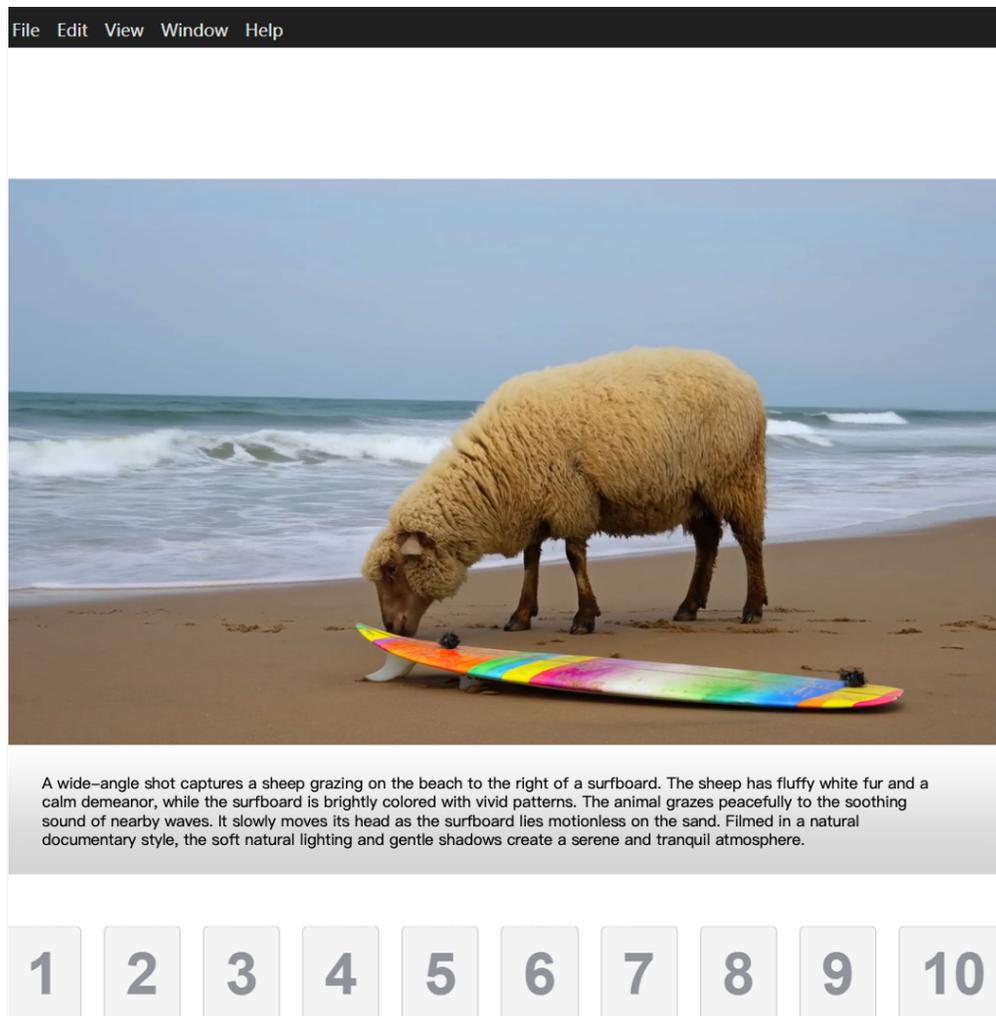


Figure 9: User Interface for the subjective experiment.