

Your One-Stop Solution for AI-Generated Video Detection

Supplementary Material

6. Content Structure in Appendix

We have organized additional important content in the Appendix due to the limited space. We present a brief outline of the content structure of the Appendix to facilitate readers to find the corresponding content, as follows:

- **Sec. 7: Further Findings and Analysis**
 - Sec. 7.1: Definition of Evaluation Metrics;
 - Sec. 7.2: Further Findings and Analyses of Finding-1;
 - Sec. 7.3: Further Findings and Analyses of Finding-2;
- **Sec. 8: Details of AIGVDBench:**
 - Sec. 8.1: Brief Introduction of Generation Models;
 - Sec. 8.2: Introduction of the Used Detection Methods;
 - Sec. 8.3: Introduction of Original Data;
 - Sec. 8.4: Details of Prompts Categorization;
 - Sec. 8.5: Balance Comparison with Prior Datasets;
 - Sec. 8.6: Discussion on Video Encoding Methods;
 - Sec. 8.7: Training details;
 - Sec. 8.8: Timeline of AIGVDBench Construction;
- **Sec. 9: Limitation;**
- **Sec. 10: More Examples of AIGVDBench:**
 - Sec. 10.1: Video Samples from the AIGVDBench;
 - Sec. 10.2: More Experimental Results;

7. Further Findings and Analysis

7.1. Definition of Evaluation Metrics.

The VBench benchmark suite systematically evaluates the performance of video generation models across six core dimensions, with each metric quantifying a distinct key characteristic of the generated content. The metrics are defined and explained below:

- **Aesthetic Quality:** This metric assesses the artistic value and visual appeal of individual video frames, encompassing compositional harmony, colour richness and coordination, photographic realism, naturalness, and overall artistic merit. It reflects human perceptual standards for visual content, analogous to criteria used in evaluating professional photography.
- **Background Consistency:** This measures the temporal coherence of the background scene throughout a video sequence. It focuses on the stability and consistency of non-subject environmental elements (e.g., landscapes, buildings) across frames, aiming to detect illogical mutations or structural distortions in the background.
- **Dynamic Degree:** This quantifies the intensity and

magnitude of motion within a video. It evaluates whether the generated content exhibits significant motion characteristics (e.g., rapidly moving objects, camera movement) appropriate to the scene context, thereby distinguishing dynamic content from static imagery and reflecting the level of visual vitality.

- **Imaging Quality:** This metric detects low-level image attributes and technical artefacts in video frames. It primarily evaluates distortions such as over-exposure, image noise, motion blur, and colour inaccuracies, focusing on objective signal-level quality rather than subjective aesthetic judgement.
- **Motion Smoothness:** This analyses the fluency and naturalness of motion in the video. It examines whether motion trajectories adhere to physical principles (e.g., uniform or accelerated motion), with a focus on identifying non-continuous phenomena such as jerky transitions, jitter, or temporal inconsistencies.
- **Subject Consistency:** This evaluates the temporal coherence of the appearance of the main subject (e.g., a person, animal, or vehicle) in the video. It requires that the subject’s identity and visual attributes remain consistent across frames, preventing implausible changes in colour, shape, or texture.

7.2. Further Findings and Analyses of Findings-1

Analysis-1.1: An analysis on the heightened difficulty of detecting closed-source models relative to open-source ones.

1. A fundamental contrast exists between open and closed source models, which is clearly reflected in their output quality. Our evaluation based on the VBench framework across six dimensions indicates that closed source models generally achieve significantly higher overall video quality than open source models (Fig. 4 and Tab. 7). This performance disparity aligns with the inherent differences between the two: open source models boast full transparency in design and parameters, often with a clear lineage of iterative refinement, whereas closed source models are opaque, not only in terms of hidden hyperparameters but also regarding potential post processing techniques. Moreover, their commercial orientation typically leads to larger scale designs with vast parameter counts, ultimately contributing to their superior generative performance.

2. A critical difference also lies in the data itself. Our

Table 7. Comparison of video generation models on Aesthetic Quality, Background Consistency, Dynamic Degree, Imaging Quality, Motion Smoothness, Subject Consistency, Frame-Level (Image & Aesthetic Average) and Final Score based on Vbench. [29]

Model Name	Aesthetic Quality	Imaging Quality	Frame-Level	Background Consistency	Dynamic Degree	Motion Smoothness	Subject Consistency	Final Score
Open-Source T2V Models								
Open-Sora	54.46	53.68	54.07	97.24	7.17	95.80	97.01	73.70
LTX	54.11	58.60	56.36	97.11	22.53	96.90	96.66	77.44
RepVideo	55.08	63.20	59.14	96.32	23.63	93.07	95.50	77.60
AnimateDiff	60.95	67.07	64.01	97.87	5.70	96.36	98.93	77.61
AccVideo	58.37	63.72	61.05	99.17	9.10	98.89	99.15	77.89
HunyuanVideo	59.31	64.04	61.68	98.92	9.67	98.37	98.87	78.03
Pyramid-Flow	58.32	63.16	60.74	97.16	18.58	96.23	96.86	78.24
VideoCrafter	64.17	63.80	63.99	97.88	11.27	95.35	97.95	78.26
Cogvideox1.5	54.98	63.39	59.19	96.91	27.25	94.80	96.74	78.92
IPOC	56.14	64.86	60.50	97.39	23.43	94.86	97.88	79.01
EasyAnimate	61.44	63.90	62.67	97.12	25.22	93.09	97.28	79.64
Wan2.1	59.52	68.03	63.78	97.14	31.28	92.18	96.62	80.87
Open-Source I2V Models								
SVD	53.52	59.57	56.55	92.34	36.52	81.10	92.22	75.50
SEINE	53.86	64.98	59.42	92.14	29.13	89.03	92.14	76.60
LTX	55.30	64.73	60.02	96.71	17.23	95.46	95.97	77.35
VideoCrafter	57.81	64.00	60.91	97.06	15.50	94.49	97.15	77.46
Pyramid-Flow	56.09	64.09	60.09	97.06	20.57	96.01	97.03	78.34
EasyAnimate	58.34	66.43	62.39	94.85	31.65	90.02	94.04	79.15
Open-Source V2V Models								
Cogvideox1.5	55.29	66.24	60.77	98.19	7.65	97.73	98.79	77.07
LTX	56.91	67.58	62.25	98.03	9.78	97.14	98.30	77.77
Closed-Source Models								
Pika	58.83	58.52	58.68	98.29	10.38	98.31	97.94	76.77
Gen3	58.30	60.19	59.25	94.72	23.50	95.27	91.62	77.02
Gen2	63.04	64.10	63.57	98.99	3.68	98.50	98.64	77.63
Vidu	60.67	64.58	62.63	97.00	19.88	94.27	96.59	78.72
OpenSora	62.99	64.95	63.97	97.75	18.20	94.19	96.92	79.09
Wan	60.65	64.68	62.67	97.97	20.88	94.87	96.78	79.24
Kling	62.76	64.62	63.69	97.28	22.38	94.02	96.28	79.52
Luma	62.97	63.33	63.15	97.84	23.10	95.51	96.38	79.84
Sora	60.57	63.38	61.98	97.05	29.25	94.06	95.45	79.96
Jimeng	68.29	63.99	66.14	98.33	21.15	92.67	97.10	80.28
Causvid 24fps	67.72	65.84	66.78	97.18	40.25	84.76	95.68	82.08
I2V AVG	<u>55.82</u>	<u>63.97</u>	59.90	95.03	<u>25.10</u>	91.02	94.76	77.40
V2V AVG	<u>56.10</u>	66.91	61.51	98.11	8.72	97.44	98.54	<u>77.42</u>
T2V AVG	58.07	63.12	<u>60.60</u>	<u>97.52</u>	<u>17.90</u>	<u>95.49</u>	<u>97.45</u>	78.10
Open-Source AVG	<u>57.20</u>	63.75	<u>60.48</u>	<u>96.83</u>	<u>19.14</u>	94.34	96.75	<u>77.82</u>
Closed-Source AVG	62.44	<u>63.47</u>	62.96	97.49	21.15	<u>94.22</u>	<u>96.31</u>	79.10



Figure 4. Comparison of video generation models on Aesthetic Quality, Background Consistency, Dynamic Degree, Imaging Quality, Motion Smoothness, Subject Consistency, Frame-Level (Image & Aesthetic Average) and Final Score based on Vbench. [29]

open source dataset was constructed with strict controls to guarantee homogeneity in semantic content. This is evidenced by CLIP-ViT-L/14 feature projections (Fig. 5), where videos form tight clusters containing outputs from all generative models. No such constraints apply to the closed source data, which is a heterogeneous collection of various tasks (primarily text to video). Thus, from a detection standpoint, the open source data constitutes a single domain benchmark, while the closed source data presents a far more challenging cross domain scenario, leading to a substantial increase in detection difficulty.

A noteworthy and interesting phenomenon observed in Fig. 6 is that, under this open set setting, the features of closed source models do not cluster by semantic content as seen in open source models, but instead exhibit a model wise hierarchical clustering pattern, which aligns with the findings reported in [56]. This observation leads us to hypothesize that when real and fake videos in a dataset lack

Table 8. Performance of TimeSformer, UniformerV2, and Effort Across Video Generation Tasks and Model Sources

Model	Close Source	Open Source			
		I2V	T2V	V2V	AVG
TimeSformer [6]	86.50	75.46	86.84	67.59	81.50
UniformerV2 [36]	87.07	81.27	85.94	73.30	83.27
Effort [88]	94.05	86.32	90.20	74.74	87.49

content alignment, it may induce such model driven clustering behavior. This type of pattern could potentially simplify the detection task, thereby explaining the anomalous performance of models on such datasets.

Abnormal Phenomenon. A notable exception to this performance degradation is observed in models such as TimeSformer [6], UniformerV2 [36], and Effort [88], which maintain stable performance on closed source data. To further investigate this anomalous behavior, we compared the

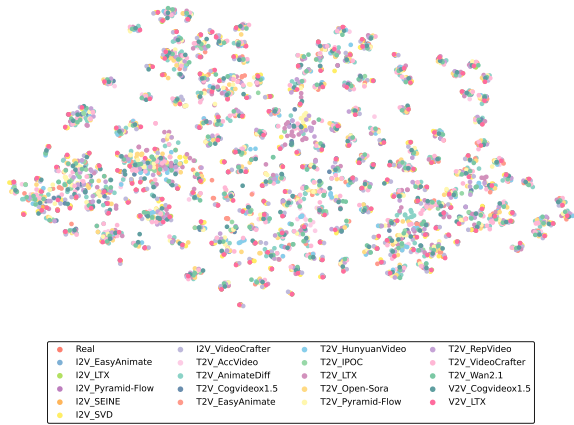


Figure 5. t-SNE [52] visualization of feature distribution for videos generated by open source models on CLIP-ViT-L/14 embeddings.

performance of these three detection models across various generative models, with a particular focus on their results on different tasks within the open source models. As shown in the Tab. 8, all three models perform significantly better on Text to Video (T2V) tasks than on Image to Video (I2V) or Video to Video (V2V) tasks, a trend that aligns with their overall performance on closed source models. This outcome is expected, since closed source video content is mainly generated through T2V, a format favored by users for its convenience, and also because the VBench benchmark primarily uses T2V based evaluation tasks. Furthermore, although closed source models generally achieve higher overall quality scores than open source models, their performance on certain dimensions, such as imaging quality, motion smoothness, and subject consistency, is comparable to that of open source T2V models. This indicates that the strong T2V performance of the three detection models is the main reason for their anomalous stability on closed source data. This conclusion is further supported by their notably high detection accuracy observed in T2V models. This analysis leads to a new question: is there an underlying pattern that determines how different detection models perform across different generation tasks?

Analysis-1.2: Further analysis on Effort and ForgeLens.

Tabs. 3 and 5 demonstrates the potential of advanced generated image detection algorithms for frame level detection in generated videos, with the Effort [88] and ForgeLens [17] models exhibiting particularly outstanding performance. A noticeable trend in generated image detection is the transition of backbone networks from architectures such as ResNet to CLIP-ViT-L/14 following the introduction of UnivFD [56].

To investigate the rationale behind this shift, we conducted a feature visualization analysis comparing the off

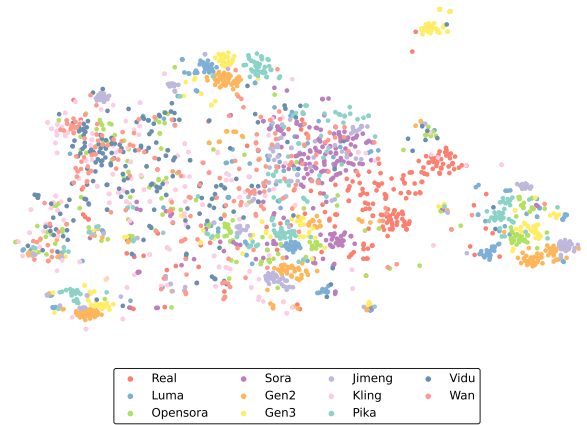


Figure 6. t-SNE [52] visualization of feature distribution for videos generated by closed source models on CLIP-ViT-L/14 embeddings.

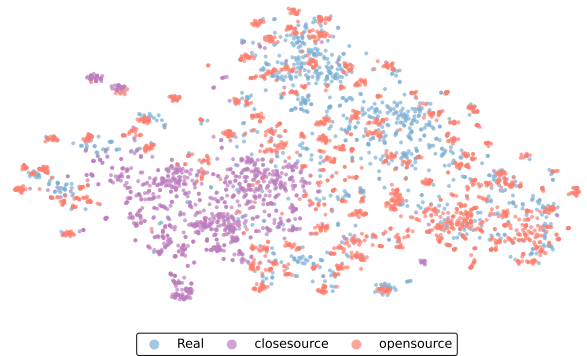


Figure 7. t-SNE [52] visualization of feature distribution from real, open-source models generated, and closed-source models generated videos on UnivFD.

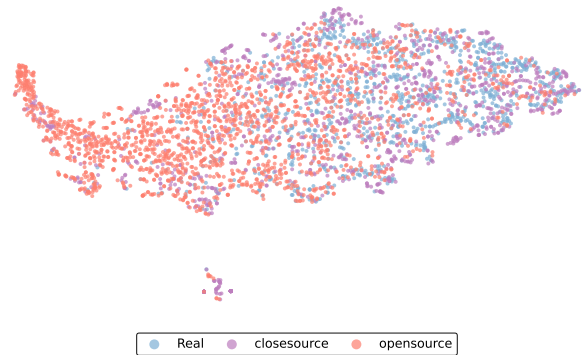


Figure 8. t-SNE [52] visualization of feature distribution from real, open-source models generated, and closed-source models generated videos on CNNSpot.

the shelf CLIP-ViT-L/14 without additional training and a ResNet model trained on the Open-Sora dataset. The anal-

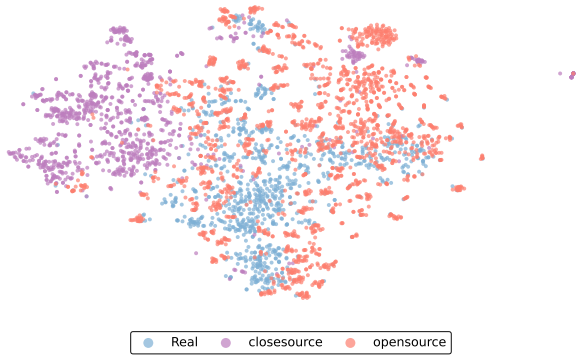


Figure 9. t-SNE [52] visualization of feature distribution from real, open-source generated, and closed-source generated videos on Effort.

ysis included features from real videos and those generated by both open source and close source models. As shown in Figs. 7 and 8 and Tab. 3, the untrained CLIP-ViT-L/14 exhibits a feature discrimination capability comparable to the trained ResNet. However, its performance is slightly inferior when handling videos generated by close source models (Tab. 5). This is corroborated by the greater separability between features of close source model generated videos and real videos observed in Fig. 8.

In summary, CLIP-ViT-L/14 shows promise as a backbone network for generated video detection, a point noted in previous work [11, 50]. Nevertheless, as observed from Fig. 7 and Tab. 5, there remains room for improvement. The fine tuning strategy proposed by Effort significantly enhances the detection performance of the base model (Fig. 9), improvements of 1.8% on open source models and 28.79% on closed source models. Furthermore, the training strategy employed by ForgeLens, which utilizes a small dataset over more training epochs (ForgeLens1 in Tabs. 3 and 5), proves more effective compared to using large datasets with fewer training epochs (ForgeLens3 in Tabs. 3 and 5), gains of 2.22% and 8.16% were achieved on open-source and closed source models, respectively.

Finding-1.1: The type of generation task has a significant impact on detector performance, and the extent of this impact varies considerably across different types of detectors.

Tab. 9 presents the performance evaluation of various detection models across multiple video generation tasks, with the following key observations:

1. **Video to Video (V2V) consistently proves to be the most challenging task to detect across all models.** This result is expected, as V2V utilizes a real video as reference, leading to generated content that generally outperforms other tasks in various metrics (as also corroborated

Table 9. Comprehensive Performance Evaluation of Video Detection Models on Text-to-Video, Image-to-Video, and Video-to-Video Generation Tasks.

Method	I2V	T2V	V2V	Closed
I3D	89.05	80.78	78.07	61.18
MViTv2	65.60	66.96	63.57	56.94
SlowFast	72.34	82.31	78.18	62.71
VideoSwin	74.97	67.93	62.01	52.47
Timesformer	75.46	86.84	67.59	86.50
TSM	83.90	74.52	67.78	56.08
Uniformer	53.35	74.12	72.21	57.02
Uniformerv2	81.27	85.94	73.30	87.07
VideoMAE	67.56	78.51	71.40	57.66
X3D	90.40	83.65	67.60	66.04
CNNSpot	89.57	84.80	78.60	70.60
D3	75.32	87.28	75.28	76.58
Effort	86.32	90.20	74.74	94.05
ForgeLens1	98.09	90.22	82.63	85.03
ForgeLens3	96.98	87.78	78.39	76.87
FreDect	72.52	69.86	58.81	55.68
Fusing	86.30	77.36	72.33	65.46
Gram-Net	83.48	80.97	69.62	65.68
NPR	75.19	66.85	61.65	52.36
UnivFD	90.30	83.97	82.20	65.26
DeCoF	76.65	85.68	81.59	72.90
DeMamba	82.09	81.87	72.47	69.43

■ Maximum ■ Middle ■ Minimum
Color coding applies only to I2V, T2V and V2V columns.

Table 10. Comparison of video generation tasks on Aesthetic Quality, Background Consistency, Dynamic Degree, Imaging Quality, Motion Smoothness, Subject Consistency, Frame-Level (Image & Aesthetic Average) and Final Score based on Vbench. [29]

Dimension	I2V	T2V	V2V
Aesthetic Quality	55.82	58.07	56.10
Imaging Quality	63.97	63.12	66.91
Frame-Level	59.90	60.60	61.51
Background Consistency	95.03	97.52	98.11
Dynamic Degree (DD)	25.10	17.90	8.72
Motion Smoothness	91.02	95.49	97.44
Subject Consistency	94.76	97.45	98.54
Video-Level(w/o DD)	93.60	96.82	98.03
Final Score	77.40	78.10	77.42

■ Best ■ Middle ■ Worst
Color coding applied per row (each dimension independently).

Table 11. Performance of detection models on videos from T2V, I2V, and V2V tasks of a single generation model.

Method	LTX Dataset			EasyAnimate		Pyramid-Flow		VideoCrafter		Cogvideox1.5	
	I2V	T2V	V2V	I2V	T2V	I2V	T2V	I2V	T2V	T2V	V2V
I3D	81.68	94.62	76.40	81.00	78.40	94.01	96.75	84.79	91.49	83.49	79.74
MViTv2	67.74	69.52	64.93	52.26	60.65	65.93	70.02	70.45	68.61	63.61	62.21
SlowFast	71.62	76.22	77.50	63.92	74.91	83.83	92.75	68.36	91.33	80.23	78.86
VideoSwin	70.88	70.08	66.52	60.43	62.99	65.63	63.88	77.10	76.33	65.64	57.49
Timesformer	70.38	85.51	65.55	66.33	84.34	80.32	96.17	65.38	96.42	84.55	69.63
TSM	70.38	80.06	66.67	79.41	77.64	85.88	87.01	84.25	94.45	68.53	68.89
Uniformer	60.05	52.92	69.35	48.46	65.95	66.36	84.92	68.78	83.73	67.57	75.06
Uniformerv2	77.59	82.60	72.65	73.94	84.35	77.11	92.95	82.28	97.14	82.95	73.94
VideoMAE	69.09	75.01	71.54	57.50	66.88	77.04	92.66	65.40	86.68	75.33	71.25
X3D	68.43	84.17	61.67	90.98	93.81	92.77	97.06	96.56	97.48	76.92	73.52
CNNSpot	81.41	96.82	71.81	88.54	88.20	97.33	99.21	73.17	91.97	89.74	85.39
D3	77.69	83.82	77.86	74.69	84.82	82.45	98.37	42.52	95.37	90.09	72.69
Effort	79.15	83.13	78.42	88.74	97.01	94.99	99.87	58.25	99.61	89.31	71.05
ForgeLens1	92.89	95.63	91.41	98.21	98.68	98.68	99.71	98.86	99.81	87.03	73.85
ForgeLens3	89.67	93.76	86.79	95.79	96.17	97.78	99.59	98.76	99.50	84.34	69.98
FreDect	61.23	70.06	52.62	59.70	67.01	69.19	73.32	77.62	81.54	71.51	65.00
Fusing	73.76	90.64	64.38	81.69	76.79	93.07	95.98	77.91	87.98	81.90	80.28
Gram-Net	71.57	92.95	59.62	77.22	81.26	88.61	95.48	77.04	93.49	84.00	79.61
NPR	67.60	80.67	58.58	69.58	65.39	77.80	85.35	59.54	70.85	75.61	64.72
UnivFD	89.91	94.70	89.38	87.39	88.73	89.57	97.39	81.12	97.55	84.69	75.02
DeCoF	81.53	87.74	81.34	70.64	78.03	83.74	96.57	45.21	90.77	90.71	81.84
DeMamba	77.59	85.28	77.46	76.37	76.59	92.84	98.90	52.02	88.01	82.17	67.48

■ Maximum ■ Middle ■ Minimum

Table 12. Comparison of Aesthetic Quality, Background Consistency, Dynamic Degree, Imaging Quality, Motion Smoothness, Subject Consistency, Frame-Level (Image & Aesthetic Average), and Final Score across video generation tasks of a single generation model. [29]

Metric	LTX			EasyAnimate		Pyramid-Flow		VideoCrafter		Cogvideox1.5	
	I2V	T2V	V2V	I2V	T2V	I2V	T2V	I2V	T2V	T2V	V2V
Aesthetic Quality	55.30	54.11	56.91	58.34	61.44	56.09	58.32	57.81	64.17	54.98	55.29
Imaging Quality	64.73	58.60	67.58	66.43	63.90	64.09	63.16	64.00	63.80	63.39	66.24
Frame-Level	60.02	56.36	62.25	62.39	62.67	60.09	60.74	60.91	63.99	59.19	60.77
Background Consistency	96.71	97.11	98.03	94.85	97.12	97.06	97.16	97.06	97.88	96.91	98.19
Dynamic Degree	17.23	22.53	9.78	31.65	25.22	20.57	18.58	15.50	11.27	27.25	7.65
Motion Smoothness	95.46	96.90	97.14	90.02	93.09	96.01	96.23	94.49	95.35	94.80	97.73
Subject Consistency	95.97	96.66	98.30	94.04	97.28	97.03	96.86	97.15	97.95	96.74	98.79
Final Score	77.35	77.44	77.77	79.15	79.64	78.34	78.24	77.46	78.26	78.92	77.07

■ Best ■ Middle ■ Worst

by Tab. 10). The lower Dynamic Degree observed in V2V outputs is primarily due to the high frame rate of the input real video, which restricts the dynamic variation range of the generated video.

2. Image-level detectors show a clear task-dependent performance trend. For image-level detectors, the influ-

ence of the generation task type on detection performance exhibits a clear and consistent pattern. Most models achieve the best detection results on Image to Video (I2V), followed by Text to Video (T2V), with V2V being the most difficult. It is noteworthy that this detection performance ranking is inversely correlated with the overall image-level qual-

ity ranking of the three tasks: V2V achieves the highest visual quality, followed by T2V, while I2V scores the lowest. Although the performances of Effort and D3 differ from the majority of models, their results align with the ranking of tasks based on the Imaging Quality metric. Therefore, it can be inferred that for image-level detectors, the generation task type does influence detection effectiveness, and this influence is negatively correlated with the image-level quality of the generated content: higher quality leads to greater detection difficulty. However, this relationship likely stems from the combined effect of multiple quality metrics rather than being dictated by any single indicator.

3. Video-level detectors exhibit more complex behavior influenced by both video and image attributes. For video-level detectors, the impact of task type is more complex, influenced not only by video-level metrics but also by image-level indicators. According to the video-level metrics in Tab. 10 (disregarding Dynamic Degree), V2V shows the best overall performance, followed by T2V, with I2V being the weakest. A straightforward interpretation is that although image-to-video (I2V) generation benefits from richer input information compared to text-to-video (T2V) generation, the reference image may, to some extent, constrain the temporal coherence and overall consistency of the generated video in certain generative models. We hypothesize that, under the condition of using a consistent generative model, the output quality would exhibit a progressive improvement from T2V to I2V, and further to video-to-video (V2V) generation. This hypothesis will be specifically validated in point 4.

4. For the same generative model, detectors typically underperform on V2V tasks compared to I2V tasks, while I2V tasks present greater challenges than T2V tasks. This trend aligns with the output quality of generative models across different tasks. As shown in Tab. 11, for the majority of detectors, the detection difficulty of V2V content generated by the same model is higher than that of I2V content, which in turn is higher than T2V content. Although a few exceptions exist, these are attributed to inherent differences among the detectors. From the perspective of generative models, the same model generally achieves higher generation quality in V2V tasks than in I2V tasks, and higher in I2V than in T2V tasks, as shown in Sec. 7.2. However, certain generative models, such as EasyAnimate, Pyramid-Flow and VideoCrafter do not strictly follow this pattern in I2V and T2V tasks. Notably, these anomalies in generation quality do not significantly affect detector performance; detectors consistently reflect an increasing order of difficulty: $V2V > I2V > T2V$. Further analysis reveals that although I2V generation quality is lower than T2V across most metrics, it outperforms T2V in terms of "imaging quality" and "dynamic extent." This suggests that these two metrics may play a critical role in detection diffi-

culty.

However, different video-level detection models exhibit considerable variation in performance across tasks, without forming a consistent pattern. While some models align with the trend indicated by video-level metrics, others show significant heterogeneity. This complexity motivates further investigation, leading to the research question raised in Sec. 7.3.

Finding-1.2: Current VLMs lack reliable capability for detecting AI-generated videos.

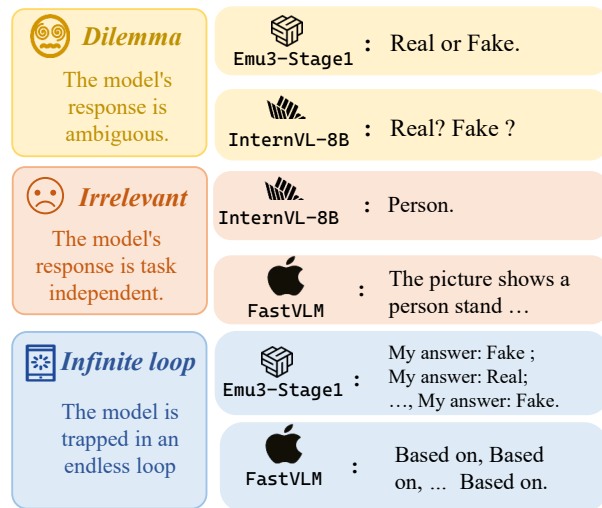


Figure 10. Anomalies in VLM for AI generated video detection.

In Sec. 8.7, we detail the experimental setup for VLMs: given an input video, the model is required to determine whether the content is AI-generated or authentic and return a binary Boolean output. For VLMs lacking multi-frame reasoning capability, we adopt a frame-by-frame evaluation strategy. Specifically, if all frames are classified as "Fake", the video is considered "Fake"; if at least one frame is judged as "Real", the video is labeled as "Real". Videos that do not elicit a valid response from the model are categorized as "No Answer".

As illustrated in Fig. 10, we classify these "No Answer" cases into three main types:

- **Dilemma:** The model produces an ambiguous response containing both "Real" and "Fake" labels in a single judgment. This differs from the case where a model lacks multi-frame reasoning ability; here, the ambiguity arises within a single inference step.
- **Irrelevant Response:** The model provides answers unrelated to the detection task, behaving as though it is performing a video question-answering task, such as describing the video content.

Table 13. ACC Performance Comparison of Vision-Language Models on Open-Source Generative Models: a video is labeled as Fake if at least one frame is classified as Fake.

Method	I2V						T2V										V2V		AVG		
	Easy Animate	LTX	Pyramid Flow	SEINE	SVD	Video Crafter	Acc Video	Animate Diff	Cogvideo x1.5	Easy Animate	Hunyuan	IPOC	LTX	Open Sora	Pyramid Flow	Rep Video	Video Crafter	Wan 2.1		Cogvideo x1.5	LTX
Emu3-Stage1 [79]	33.19	33.53	33.01	32.24	28.67	32.78	31.01	40.57	31.78	32.29	30.03	31.25	33.64	27.63	30.62	31.64	42.84	31.92	34.63	35.29	32.93
FastVLM-7B [75]	46.61	47.74	47.38	47.68	50.63	49.69	47.61	49.03	48.30	46.85	47.58	48.63	48.78	48.58	46.89	48.91	51.14	46.72	48.21	47.54	48.22
DeepseekVL2-S [82]	49.02	51.35	48.60	49.37	54.82	48.90	48.38	48.27	48.62	48.42	48.37	48.47	49.35	49.08	49.45	48.32	48.43	48.37	48.62	49.14	49.14
LLaVA-v1.5-7B [44]	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00
Kimi-VL-A3B [72]	50.27	52.17	50.08	50.83	55.37	51.57	50.10	51.03	50.47	50.07	49.98	50.33	50.83	49.80	51.45	51.55	50.87	50.00	50.07	50.03	50.84
Qwen2.5-VL-3B [5]	49.05	51.50	48.77	53.00	60.03	61.10	49.10	56.47	50.27	45.52	47.65	50.80	51.22	45.87	45.97	52.37	58.43	48.87	50.40	49.03	51.27
Deepseek-VL-7B [47]	51.03	55.10	50.58	53.13	59.30	52.50	50.15	49.80	51.82	50.13	50.13	50.50	54.42	50.50	50.72	53.37	50.05	49.95	50.78	50.30	51.71
DeepseekVL2 [82]	52.18	56.42	51.28	58.43	67.07	59.77	49.22	50.08	53.72	49.17	48.98	52.22	57.80	49.88	52.83	55.30	50.58	49.23	50.20	50.43	53.24
Qwen2.5-VL-32B [5]	49.68	51.92	49.50	51.58	61.30	59.98	51.73	62.10	53.67	48.90	49.83	53.87	50.73	49.93	50.78	54.05	64.43	50.05	52.00	50.25	53.32
Qwen2.5-VL-7B [5]	52.45	54.75	51.98	55.22	61.72	59.57	53.67	59.08	55.32	51.25	52.68	55.08	54.82	53.72	53.10	58.28	63.35	53.73	53.32	52.00	55.25
InternVL-8B [101]	50.33	55.39	51.28	58.18	65.45	64.44	52.11	63.59	55.57	46.73	51.89	55.98	58.12	53.72	52.98	57.65	67.32	52.85	52.34	50.99	55.85

Table 14. ACC Performance Comparison of Vision-Language Models on Closed-Source Generative Models: a video is labeled as Fake if at least one frame is classified as Fake.

Method	Luma	OpenSora	Sora	Causvid	Gen2	Gen3	Jimeng	Kling	Pika	Vidu	Wan	AVG
Emu3-Stage1 [79]	36.78	38.91	37.31	41.84	40.16	34.89	42.76	39.69	37.28	37.75	36.18	38.50
DeepseekVL2-S [82]	48.63	48.58	48.40	48.28	48.18	49.25	48.23	48.58	48.78	48.83	48.45	48.56
LLaVA-v1.5-7B [44]	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00
Deepseek-VL-7B [47]	50.90	49.95	50.40	50.33	49.63	51.40	49.75	51.03	50.83	50.40	50.33	50.45
Kimi-VL-A3B [72]	53.60	52.55	53.78	55.23	49.95	58.23	50.30	59.90	56.10	54.70	53.73	54.37
FastVLM-7B [75]	51.85	52.88	50.85	53.70	48.50	52.68	51.44	60.16	57.00	52.84	53.48	53.21
Qwen2.5-VL-7B [5]	59.83	57.83	56.00	60.35	50.68	60.45	55.73	70.23	67.55	59.23	59.03	59.72
Qwen2.5-VL-32B [5]	59.85	60.83	55.60	57.93	50.65	58.40	61.63	67.75	70.38	58.38	59.90	60.12
InternVL3-8B [101]	62.78	58.43	53.18	60.08	47.85	59.05	56.18	62.46	68.59	52.38	56.39	57.94
Qwen2.5-VL-3B [5]	64.25	61.33	64.73	63.78	56.90	68.10	61.95	71.68	75.53	63.08	66.00	65.21
DeepseekVL2 [82]	87.10	87.60	84.00	85.43	82.20	86.10	83.83	85.85	91.75	84.18	84.18	85.65

- Infinite Loop: The model enters an infinite repetition cycle, generating either meaningless strings of characters or continuously revising its own responses.

Although such outputs may occur occasionally under normal circumstances, we observe that "No Answer" responses constitute a considerable proportion of the total outputs for many VLMs. To ensure a fair evaluation, these invalid responses are assigned a 50% accuracy rate in our calculations.

Furthermore, as indicated in As shown in Tabs. 3 and 6, even after excluding "No Answer" cases, the performance of VLMs that provide normal responses remains inadequate for the generated video detection task. It should be noted that the evaluation criteria applied to VLMs without temporal reasoning capabilities may be particularly stringent: a model must correctly identify every frame in a video as fake to classify the video as "Fake". To examine the potential influence of this experimental design, we recalculated the accuracy using two alternative strategies in Tabs. 13 to 16: majority voting, and a more lenient approach where a video is labeled "Fake" if at least one frame is detected as fake. However, even under these adjusted evaluation protocols, VLM performance remains insufficient for reliably performing the generated video detection task. Moreover, for VLMs within the same architecture family, increasing

the number of model parameters does not lead to significant performance gains.

Analysis-1.3: An exploration of whether the performance superiority of DeepSeek-VL2 signifies a genuine ability to discern authenticity.

To further investigate the strong performance of DeepSeek-VL2 on closed-source models, we conducted two experiments using the first 200 generated videos from the Luma test set in AIGVDbench and their corresponding real videos:

- Synthetic Video Attribution (SVA): The model was asked to provide open-ended responses to determine video authenticity and explain its reasoning.
- Metric Attribution Selection (MAS): The model was required to select potential factors from six evaluation metrics for generated videos that influenced its judgment on video authenticity.

The prompts and examples for these two experiments are shown in Figs. 12 and 13. Notably, despite the prompts being designed without apparent bias, the model's judgments in these two experiments diverged significantly. As shown in Tab. 17, in Experiment 1, the model classified the vast majority of videos as real, while Experiment 2 demonstrated the opposite tendency. Figs. 12 and 13 revealed that

Table 15. ACC Performance Comparison of Vision-Language Models on Open-Source Generative Models under the majority voting setting: the video label is determined by the majority classification of its individual frames.

Method	I2V						T2V										V2V		AVG		
	Easy Animate	LTX	Pyramid Flow	SEINE	SVD	Video Crafter	Acc Video	Animate Diff	Cogvideo x1.5	Easy Animate	Hunyuan	IPOC	LTX	Open Sora	Pyramid Flow	Rep Video	Video Crafter	Wan 2.1		Cogvideo x1.5	LTX
Emu3-Stage1 [79]	33.23	33.57	33.04	32.28	28.70	32.82	31.04	40.60	31.82	32.33	30.07	31.28	33.68	27.66	30.65	31.68	42.88	31.95	34.66	35.33	32.96
FastVLM-7B [75]	46.56	47.74	47.38	47.59	50.53	49.59	47.58	48.98	48.15	46.73	47.57	48.50	48.71	48.61	46.88	48.73	51.04	46.57	48.14	47.54	48.15
DeepseekVL2-S [82]	49.20	49.58	49.23	49.27	50.92	49.35	49.18	49.17	49.22	49.17	49.15	49.22	49.37	49.15	49.50	49.22	49.15	49.15	49.20	49.22	49.33
LLaVA-v1.5-7B [44]	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00
Kimi-VL-A3B [72]	49.98	50.08	50.00	50.20	52.02	50.88	50.22	50.87	50.17	50.05	50.12	50.23	50.28	49.85	50.92	50.50	50.58	49.93	50.08	50.03	50.35
Deepseek-VL-7B [47]	50.05	51.30	50.20	50.52	52.93	50.78	50.20	49.90	50.68	50.08	50.10	50.12	51.65	50.07	50.32	50.75	50.05	49.95	50.42	50.20	50.51
Qwen2.5-VL-3B [5]	48.98	50.62	49.63	50.08	53.78	56.33	49.87	54.88	50.30	48.02	49.32	50.47	50.78	47.92	48.30	51.22	54.48	49.17	50.78	49.83	50.74
Qwen2.5-VL-32B [5]	49.43	50.25	49.80	49.70	53.60	55.48	51.28	58.07	52.32	49.65	50.23	52.27	50.18	49.60	50.87	52.67	58.72	49.63	51.57	50.18	51.78
DeepseekVL2 [82]	50.50	53.07	51.25	53.60	59.63	57.15	50.08	50.95	52.08	49.08	49.73	51.15	54.93	50.15	52.80	53.33	50.93	49.37	50.90	50.97	52.08
Qwen2.5-VL-7B [5]	50.70	51.58	51.28	51.78	55.48	54.67	52.52	56.40	52.80	50.58	51.88	53.03	52.35	51.93	51.92	54.90	57.75	51.77	51.80	51.18	52.82
InternVL3-8B [101]	49.02	54.43	51.75	55.18	62.70	63.56	54.14	63.63	56.07	47.67	53.43	56.27	57.40	53.82	53.65	57.70	67.83	52.80	53.49	52.11	55.83

Table 16. ACC Performance Comparison of Vision-Language Models on Closed-Source Generative Models under the strict evaluation setting: the video label is determined by the majority classification of its individual frames.

Method	Luma	OpenSora	Sora	Causvid	Gen2	Gen3	Jimeng	Kling	Pika	Vidu	Wan	AVG
Emu3-Stage1 [79]	36.83	38.96	37.36	41.89	40.21	34.94	42.81	39.74	37.33	37.80	36.23	38.55
DeepseekVL2-S [82]	49.05	49.00	49.03	49.00	49.00	49.05	49.00	49.00	49.18	49.08	49.08	49.04
Deepseek-VL-7B [47]	50.28	49.90	50.18	50.08	49.90	50.33	49.93	50.15	50.35	50.05	50.08	50.11
LLaVA-v1.5-7B [44]	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00
Kimi-VL-A3B [72]	52.45	50.95	52.28	53.05	49.68	54.30	50.03	56.45	53.65	52.75	52.40	52.54
FastVLM-7B [75]	51.83	52.88	50.85	53.73	48.58	52.63	51.49	60.01	57.00	52.89	53.53	53.22
Qwen2.5-VL-7B [5]	56.00	54.58	53.18	56.05	49.88	54.90	52.43	66.53	61.73	56.83	56.08	56.20
Qwen2.5-VL-32B [5]	56.93	58.35	53.80	55.50	50.20	55.48	58.50	65.28	66.20	57.03	58.05	57.75
InternVL3-8B [101]	60.85	57.83	52.40	58.88	48.45	56.93	55.73	63.31	68.86	52.83	55.86	57.45
Qwen2.5-VL-3B [5]	58.60	57.43	59.58	57.63	55.33	62.38	56.70	70.10	70.45	59.98	61.68	60.89
DeepseekVL2 [82]	83.13	84.13	77.70	79.28	76.93	80.73	78.58	82.20	91.35	77.85	79.70	81.05

Table 17. Performance Comparison of DeepSeek-VL2 Across Different Experimental Tasks.

Model	F_ACC	R_ACC	ACC
DeepseekVL2	58.00	95.50	76.75
DeepseekVL2+SVA	0.00	99.50	49.75
DeepseekVL2+MAS	99.50	0.00	49.75

in the open-ended responses of Experiment 1, the model’s justifications were notably inconsistent. Whether the classification was correct or not, its explanations ranged from reasonable to illogical. In Experiment 2, however, statistical analysis of frame-level and video-level selection results in Fig. 11 revealed certain patterns in the model’s behavior. For instance, in the imaging quality dimension, the model exhibited a preference of Real Luma Pika. Combined with its judgment results, this suggests that DeepSeek-VL2 considers the image quality of real videos to be inferior to Luma, which in turn is inferior to Pika. However, according to our previous evaluations of generative models, Luma’s generation quality in AIGVDbench should be superior to Pika.

Based on these findings, we conclude that although

DeepSeek-VL2’s performance on closed-source models is encouraging, it remains inadequate for the task of AI-generated video detection. This conclusion is supported by both its performance on open-source models and the two experiments designed to probe its explainability. Furthermore, the reason for its strong performance on closed-source models remains unclear. Based on comparative analysis of the experiments, we hypothesize that a potential reason could be that during the training of the closed-source model, DeepSeek-VL2 might have been exposed to simple classification tasks or training data for other tasks that included a certain proportion of videos generated by closed-source models, leading to specific recognition biases toward such videos.

7.3. Further Findings and Analyses of Findings-2

Analysis-2.1: Analysis of variations in the detectability of generative models and their relationship to detector type.

EXP A. The results in Fig. 3 indicate that the performance of the three types of detectors does not exhibit a simple linear decline as the overall generation quality of the models improves. To further investigate the underlying patterns, we extended our analysis in Figs. 14 to 16 to exam-

Table 18. Impact of Guidance Scale and Sampling Steps on Detectors.

Detector Type	Model	Guidance Scale							Sampling Steps				
		1.0	2.0	4.0	6.0	8.0	10.0	12.0	30	40	50	60	70
DecoF	AccVideo	43.98	70.88	89.71	93.47	95.41	96.75	97.53	89.23	86.86	93.47	87.24	88.54
	AnimateDiff	85.47	87.39	85.77	83.85	80.22	78.68	77.45	56.37	59.88	83.85	77.02	72.74
	Cogvideox1.5	56.02	69.07	82.29	88.22	91.49	92.43	94.29	97.06	94.41	88.22	94.36	91.93
	EasyAnimate	77.27	92.78	96.85	98.90	97.80	98.45	98.66	85.22	87.54	98.90	91.67	92.24
	HunyuanVideo	42.79	68.91	88.60	92.14	94.64	95.01	95.88	87.03	82.24	92.14	83.12	82.78
	IPOC	35.17	49.03	72.22	82.30	88.21	91.28	94.54	95.99	93.89	82.30	91.77	90.37
	LTX	82.82	74.33	79.71	84.24	85.93	86.57	89.28	97.38	94.71	84.24	94.57	92.62
	Open-Sora	28.27	47.59	67.50	73.96	79.80	81.02	82.91	87.82	84.25	73.96	81.40	74.53
	Pyramid-Flow	32.08	58.78	80.40	88.23	91.89	93.28	95.64	97.41	94.59	88.23	93.47	92.68
	RepVideo	41.89	48.94	69.62	81.91	88.51	90.93	94.53	96.81	93.04	81.91	93.48	91.67
	VideoCrafter	66.36	71.94	77.25	80.46	80.57	82.72	85.78	50.39	53.68	80.46	67.68	71.87
	Wan2.1	45.05	65.85	82.90	86.18	88.55	88.69	90.14	91.03	86.71	86.18	87.06	83.39
UnivFD	AccVideo	56.50	69.46	85.74	88.69	90.17	92.04	93.39	82.37	81.42	88.69	83.28	83.09
	AnimateDiff	92.74	90.49	91.15	90.87	89.87	89.71	89.02	76.27	78.96	90.87	88.88	88.73
	Cogvideox1.5	73.42	76.82	84.54	88.14	90.40	91.01	92.85	96.92	94.71	88.14	95.42	92.56
	EasyAnimate	95.99	98.68	99.42	99.75	99.68	99.88	99.80	92.60	95.11	99.75	96.35	97.54
	HunyuanVideo	49.68	65.71	82.98	85.85	88.35	89.34	90.32	78.45	76.43	85.85	79.47	78.42
	IPOC	50.73	55.48	72.31	80.58	86.17	88.41	91.07	96.21	94.10	80.58	93.84	91.34
	LTX	96.68	87.23	84.79	84.97	86.86	86.66	89.73	97.50	95.58	84.97	96.57	95.16
	Open-Sora	64.33	74.08	83.49	86.73	89.29	88.78	90.42	95.48	93.83	86.73	95.58	91.97
	Pyramid-Flow	69.32	81.52	91.10	94.43	95.30	96.14	97.43	98.50	97.83	94.43	97.99	97.50
	RepVideo	66.19	62.25	72.21	80.38	85.06	88.27	93.19	96.79	93.53	80.38	94.92	92.58
	VideoCrafter	88.24	89.43	91.16	94.10	94.07	95.38	96.60	73.49	77.67	94.10	85.74	90.99
	Wan2.1	43.77	58.53	76.07	77.29	78.28	79.60	79.18	86.12	81.06	77.29	82.62	78.95
I3D	AccVideo	39.54	54.15	76.70	83.42	86.70	89.67	91.34	71.97	76.99	83.42	74.36	67.67
	AnimateDiff	69.20	70.95	70.16	64.59	60.29	55.97	52.07	26.92	34.62	64.59	42.48	40.14
	Cogvideox1.5	55.68	68.05	77.14	84.66	86.66	88.17	89.67	72.01	77.44	84.66	70.68	73.87
	EasyAnimate	100.00	100.00	100.00	100.00	100.00	100.00	100.00	99.48	100.00	100.00	99.29	98.58
	HunyuanVideo	46.15	64.96	81.80	86.25	91.68	93.57	94.88	75.47	78.36	86.25	77.88	75.51
	IPOC	47.18	68.74	86.54	92.42	96.01	96.83	98.08	71.87	82.57	92.42	73.39	78.98
	LTX	32.05	48.50	62.45	68.13	70.93	71.99	76.81	46.48	52.51	68.13	48.78	53.19
	Open-Sora	36.68	58.42	74.56	77.22	83.46	81.97	85.77	70.80	81.14	77.22	71.33	73.96
	Pyramid-Flow	63.05	83.29	94.91	97.27	98.75	99.22	99.61	91.34	96.40	97.27	90.00	94.37
	RepVideo	58.18	75.84	92.23	96.83	98.58	98.99	99.47	97.05	97.49	96.83	91.28	93.93
	VideoCrafter	38.38	49.73	72.57	77.96	82.64	84.35	83.14	43.55	52.15	77.96	54.23	59.96
	Wan2.1	29.23	43.01	62.87	71.71	78.49	81.10	83.38	49.34	48.98	71.71	49.10	54.26

ine the correlation between detector performance on different generative models and multiple generation metrics. Our analysis reveals that detection difficulty is not directly determined by the “overall generation quality” of the model. Neither in T2V nor I2V tasks did we observe the expected clear relationship between overall quality and detection difficulty—particularly for T2V models, where almost no discernible pattern exists. While a slight trend appears in I2V tasks, we postulate that this may be a statistical artifact due to the limited number of available models.

What proves to be more indicative are the correlations between certain specific dimensions of generation quality and detection performance, as detailed below:

- The performance of all three types of detectors shows

a significant correlation with the **imaging quality** and **frame-level quality** of the generation models.

- The performance of video classification models is also correlated with the **dynamic degree** of the generation models.
- Both AI-generated video detection models and AI-generated image detection models exhibit correlations with **aesthetic quality**.
- In I2V tasks, the performance of AI-generated video detection models and video classification models is further correlated with **background consistency** and **subject consistency**.

Based on these phenomena, a seemingly plausible inference would be that the higher the quality of a generated

Table 19. Investigation of Video Quality Assessment: Impact of Guidance Scale and Sampling Steps.

Parameter	Value	Quality Metrics						Quality Score	Final Score
		Aesthetic Quality	Background Consistency	Dynamic Degree	Imaging Quality	Motion Smoothness	Subject Consistency		
Guidance Scale	1.0	57.35	96.12	12.50	61.84	95.72	96.02	76.28	76.28
	2.0	61.18	97.07	18.50	64.66	95.59	97.31	78.97	78.97
	4.0	61.02	97.47	23.50	63.35	95.33	97.70	79.70	79.70
	6.0	60.88	97.49	20.75	63.12	94.75	97.65	79.03	79.03
	8.0	60.79	97.33	22.75	61.88	94.08	97.45	78.96	78.96
	10.0	60.22	97.13	23.75	60.55	93.59	97.06	78.60	78.60
	12.0	59.60	97.09	23.25	59.60	93.67	97.11	78.24	78.24
Sampling Steps	30	45.79	94.20	14.75	53.52	79.81	91.59	69.03	69.03
	40	47.99	95.16	14.75	55.66	82.37	93.10	70.73	70.73
	50	60.88	97.49	20.75	63.12	94.75	97.65	79.03	79.03
	60	49.50	93.97	9.75	47.79	88.99	92.46	69.54	69.54
	70	54.96	93.93	23.00	59.75	88.35	93.15	75.12	75.12

Table 20. Performance Comparison of Detection Models Across Categories.

Detector	LTX					Pyramid-Flow					Wan2.1				
	Animals	Food/Bev.	Illust.	Vehicles	Range.	Animals	Food/Bev.	Illust.	Vehicles	Range.	Animals	Food/Bev.	Illust.	Vehicles	Range.
I3D	91.70	97.33	94.74	92.72	5.63	96.36	99.21	97.65	94.34	4.87	54.34	59.41	68.91	66.03	14.57
DeCoF	86.91	90.40	85.41	83.51	6.89	98.36	96.84	97.85	96.68	1.68	64.30	77.56	74.72	70.70	13.26
CNNSpot	99.28	99.83	99.36	99.39	0.55	94.75	97.18	96.69	95.89	2.43	59.41	73.10	75.65	69.16	16.24
UnivFD	98.83	98.55	97.78	96.92	1.90	94.46	95.45	92.92	93.90	2.53	53.16	56.41	65.56	57.75	12.40

video in certain specific dimensions, the more difficult it is to detect. However, regarding this viewpoint, two critical limitations must be emphasized: 1. Correlation does not imply causation. The observed associations cannot be directly interpreted as “higher quality leads to greater detection difficulty.” 2. More importantly, further analysis demonstrates that this inference is directly refuted.

EXP B. To systematically validate the relationship between the quality of generative models and detection difficulty, we conducted a comparative analysis of detectors trained on different generative models (Figs. 17 to 19). The experimental results reveal a complex phenomenon: the quality of generated content and detection difficulty do not exhibit the intuitively assumed positive correlation. Specifically, the earlier conclusion drawn from studies based on Open-Sora—that “higher quality leads to greater detection difficulty”—lacks generalizability. In tests involving other generative models, we observed the opposite trend: for detectors trained on high-quality generative models, the detection difficulty often decreases as certain generation metrics improve. Based on these findings, we propose that the key factor determining detection performance is not the absolute quality of the generative model, but rather the correlation between the generative pattern and the detector’s training data. When the generative pattern of a test sample falls

within the distribution of the training data (high correlation), the detector performs well; conversely, under out-of-distribution generalization scenarios (low correlation), detection efficacy decreases significantly. This naturally leads to a further question: within a single generative model, is there a correlation between the quality of its outputs and the difficulty of detecting them? This question is systematically examined in Analysis-2.4.

Analysis-2.2: Key characteristics of generative models for producing high-quality training data and their relationship to detector type.

Figs. 20b, 21b and 22b presents the cross-evaluation results of three detectors, each trained specifically on one of 12 T2V models and evaluated on all of them. This comprehensive assessment reinforces our prior conclusion: a superior generative model does not yield a more effective detector. Notably, detectors trained on certain models, like Wan 2.1, even exhibit performance degradation. Furthermore, higher evaluation metrics of a generative model do not guarantee improved detection performance.

To further validate the conclusion presented in Analysis-2.1, we conducted an additional experiment. Figs. 20a, 21a and 22a illustrates the relationship between the overall detection performance of detectors (each trained on one of 12

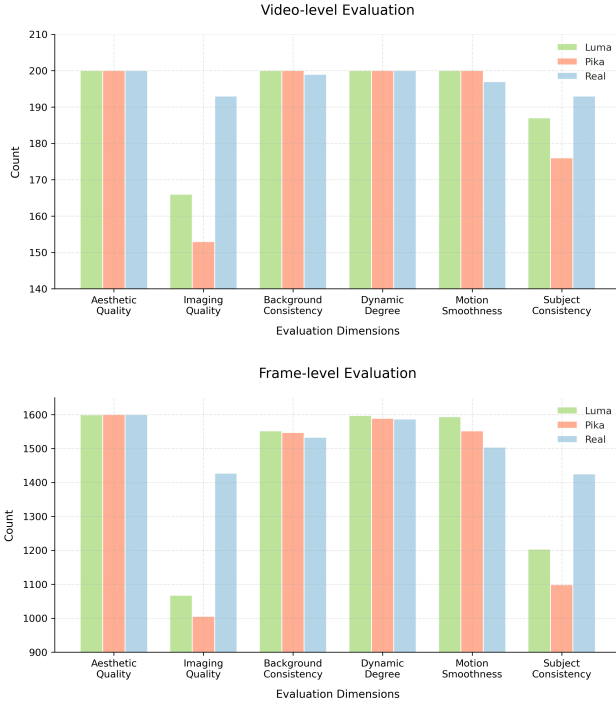


Figure 11. DeepSeek-VL-2’s Preference for Evaluation Metrics in AI-Generated Video Detection.

distinct generative models) on a specific generator and that generator’s own performance. This experiment provides further evidence that a higher-quality generative model does not necessarily lead to a greater challenge for detection.

Analysis-2.3: The effects of sampling steps on detector performance.

For a given generative model, sampling steps and guidance scale are two key parameters that influence output quality. In Analysis-2.3 and 2.4, we investigate how these factors affect detector performance and analyze the relationship between such effects and generation quality. All experiments in Tabs. 18 and 19 are conducted on a newly constructed dataset. We use EasyAnimate as the generative model, with sampling steps set to 30, 40, 50, 60, and 70, and guidance scale set to 1.0, 2.0, 4.0, 6.0, 8.0, 10.0, and 12.0. For each parameter combination, we selected the first 200 samples from the AIGVDBench test set, generating 200 synthetic videos and 200 real videos respectively, resulting in a new test set of 400 videos.

We first examine the effect of sampling steps. As shown in Tab. 19, as the number of sampling steps increases, the quality of the generated videos across both individual dimensions and overall performance first improves and then deteriorates. This pattern is expected, since EasyAnimate recommends 50 sampling steps. Fig. 23 shows a set of videos generated with different step counts. When the step count deviates from 50, the videos exhibit noticeable arti-

facts: at 40 steps, the frames are dominated by noise, while at 60 steps, the content does not match the text prompt.

In terms of detector performance, we observe two distinct trends. Detectors trained on data with stronger attributes, such as Accvideo, AnimateDiff, EasyAnimate, HunyuanVideo, and VideoCrafter, which emphasize image quality and background & subject consistency tend to achieve better detection performance as step count increases, up to a point, after which performance drops. In contrast, detectors trained on lower-quality generated data often show the opposite trend: performance first declines and then improves. This reveals two important insights:

First, the quality of training data is not simply a linear combination of individual scores, but involves a more complex interplay of factors such as visual quality and temporal coherence. Second, the results reaffirm our earlier conclusion that higher-quality generation models are not necessarily harder to detect; rather, detection performance fundamentally depends on the similarity between the test-time generative model and those seen during training. We also note the impact of detector architecture. For example, frame-based detectors such as UnivFD and DeCoF, which rely on CLIP-ViT-L/14, are more sensitive to image quality and exhibit highly consistent behavior across models. In contrast, video-level detectors such as I3D consistently show a rise-then-fall performance trend. Finally, we observe that results obtained on the 400-sample dataset (Fig. 23 for Opensora at 50 steps, and at guidance scale 6.0) differ significantly from the results reported in the main body for the same generative models. This further highlights the importance of benchmark scale in obtaining reliable evaluation outcomes.

Analysis-2.4: The effects of guidance scale on detector performance.

In contrast to the impact of sampling steps, as the guidance scale increases from 1.0 to 12.0, the generated videos exhibit divergent patterns of change across different quality dimensions. While most metrics follow the expected pattern of initial improvement followed by decline, motion smoothness demonstrates a consistent deterioration throughout this range. Regarding detector performance, with the increase in guidance scale, all detectors (except those trained on AnimateDiff) show a steady improvement in detection capability.

In stark contrast, detectors based on AnimateDiff exhibit progressive performance degradation as the guidance scale grows. The underlying reasons for this phenomenon remain unclear, particularly given that motion smoothness is fundamentally a video-level metric, while the frame-level detector UnivFD still conforms to the overall trend observed in video-based detectors. One plausible hypothesis is that higher guidance scales enhance the model’s adherence to textual prompts, resulting in greater content consis-

[VLM-based Synthetic Video Attribution]



DeepseekVL2: You are a forensic deepfake detector:

"**CRITICAL INSTRUCTION:** Begin your answer with either 'REAL' or 'FAKE', then immediately list short explanation"

"**Task:** Decide whether the frame is REAL or FAKE and briefly explain your decision"

"Respond exactly in the format 'REAL - short reason' or 'FAKE - short reason'."

"**Classification:**".

Example 1 : A real video is accurately classified as such, supported by a well-reasoned explanation.



Reason: The image shows a clear and detailed aerial view of a city with consistent lighting and realistic proportions of buildings and water.

Example 2 : A real video is accurately classified as such, but supported by an unreasonable explanation.



Reason: The image depicts a realistic landscape with natural colors and textures, consistent with a real-world scene.

Example 3 : A fake video is incorrectly classified as such, but supported by a well-reasoned explanation.



Reason: The image shows a realistic depiction of a yeti with detailed fur texture and snow particles, consistent with high-quality CGI or a well-executed prosthetic makeup.

Example 4 : A fake video is incorrectly classified as such, but supported by an unreasonable explanation.



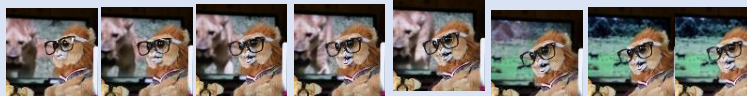
Reason: The painting style and brushstrokes are consistent with impressionist art, and the Eiffel Tower is depicted in a recognizable manner.

Example 5 : A real video is incorrectly classified as such, but supported by a well-reasoned explanation.



Reason: The image appears to be a computer-generated animation rather than a real photograph.

Example 6 : A real video is incorrectly classified as such, but supported by an unreasonable explanation.



Reason: The image appears to be digitally manipulated, as the lion's head is superimposed onto a different body, and the background is a computer-generated scene.

Example 7 : A fake video is accurately classified as such, supported by a well-reasoned explanation.



Reason: The image appears to be digitally manipulated, with elements such as the umbrellas and the background architecture having an unnatural, cartoonish quality that is not consistent with a real photograph.

Example 8 : A fake video is accurately classified as such, but supported by an unreasonable explanation.



Reason: The image appears to be digitally manipulated, as the rover's design and the Martian landscape do not match the known characteristics of actual Mars rovers or the planet's surface.

Figure 12. Prompt and Model Responses for Synthetic Video Attribution.

[VLM-based Metric Attribution Selection]



DeepseekVL2: You are a forensic deepfake detector.

"**CRITICAL INSTRUCTION:** Begin your answer with either 'REAL' or 'FAKE', then immediately list the evidence aspects in square brackets using the exact option names before giving a short explanation."

"**Task:** Analyze this video frame for deepfake detection. Determine if the content is real or AI-generated / manipulated. "

"**Focus on:** artifacts, lighting inconsistencies, texture glitches, boundary errors, unnatural movements."

"**Evidence aspect options (choose all that apply, use exact spelling, comma-separated inside the brackets):** *Aesthetic Quality, Imaging Quality, Background Consistency, Dynamic Degree, Motion Smoothness, Subject, Consistency* "

"Always include at least one listed aspect when the label is 'FAKE'. For 'REAL', include supporting aspects; if none apply, output 'REAL [Consistency]' as the default justification"

"**Reasoning:** After the bracketed aspect list, concisely describe the key visual evidence supporting your decision. Structure the sentence as LABEL [Aspect1, Aspect2] - brief explanation. "

"**Classification:**"

Example 1 : A *real* video is *incorrectly* classified as such.



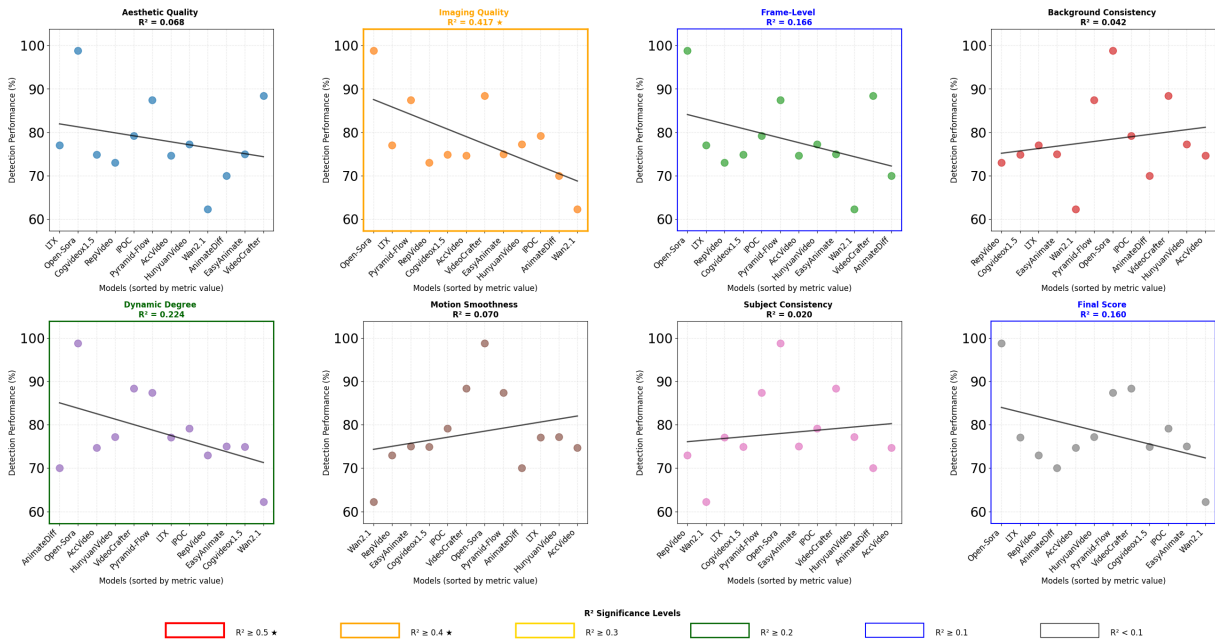
Evidence aspect options: *Aesthetic Quality, Imaging Quality, Background Consistency, Dynamic Degree, Motion Smoothness, Subject, Consistency*

Example 2 : A *fake* video is *accurately* classified as such

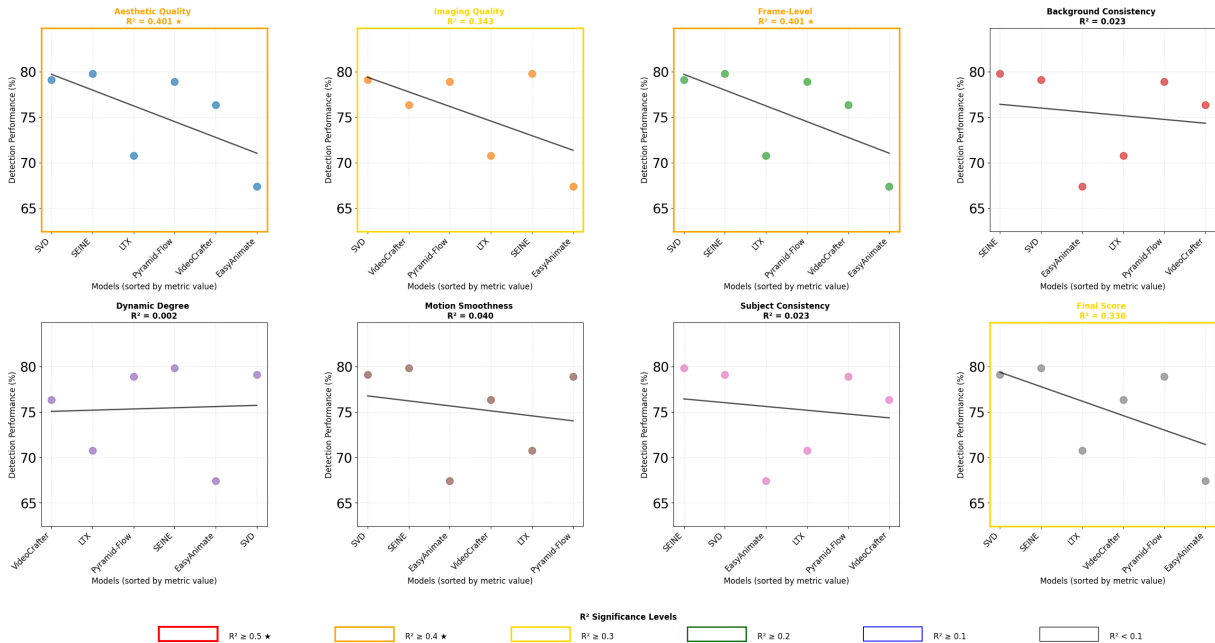


Evidence aspect options: *Aesthetic Quality, Background Consistency, Dynamic Degree, Motion Smoothness*

Figure 13. Prompt and Model Responses for Metric Attribution Selection.

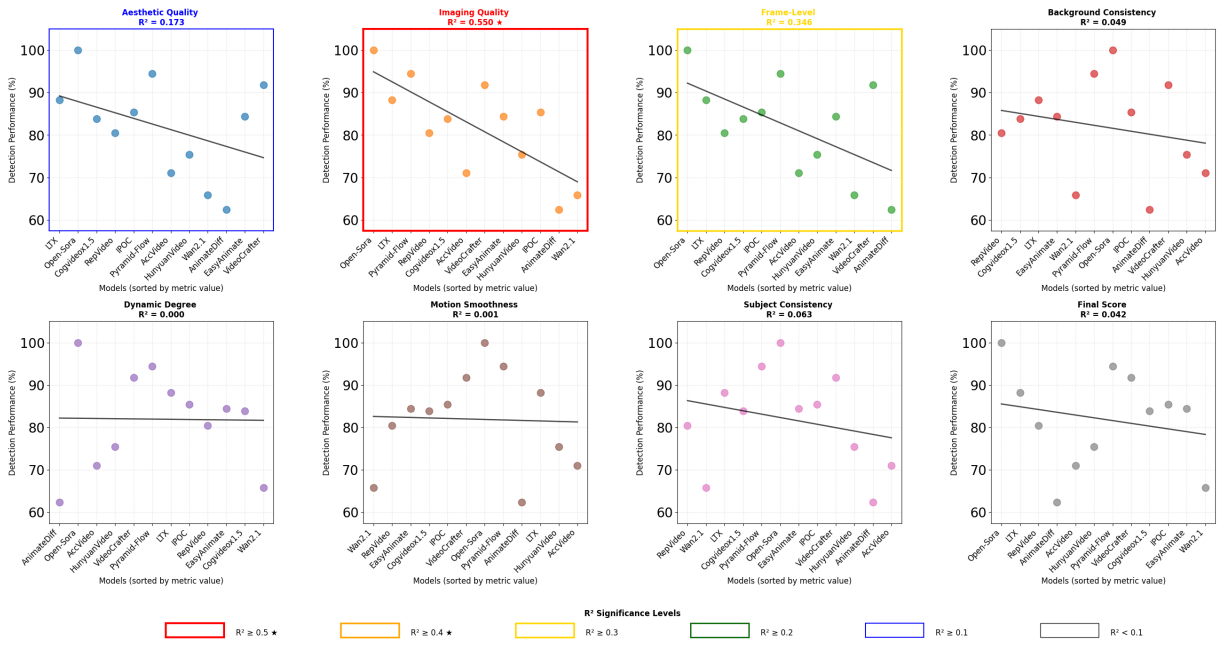


(a) Text-to-Video (T2V) Task.

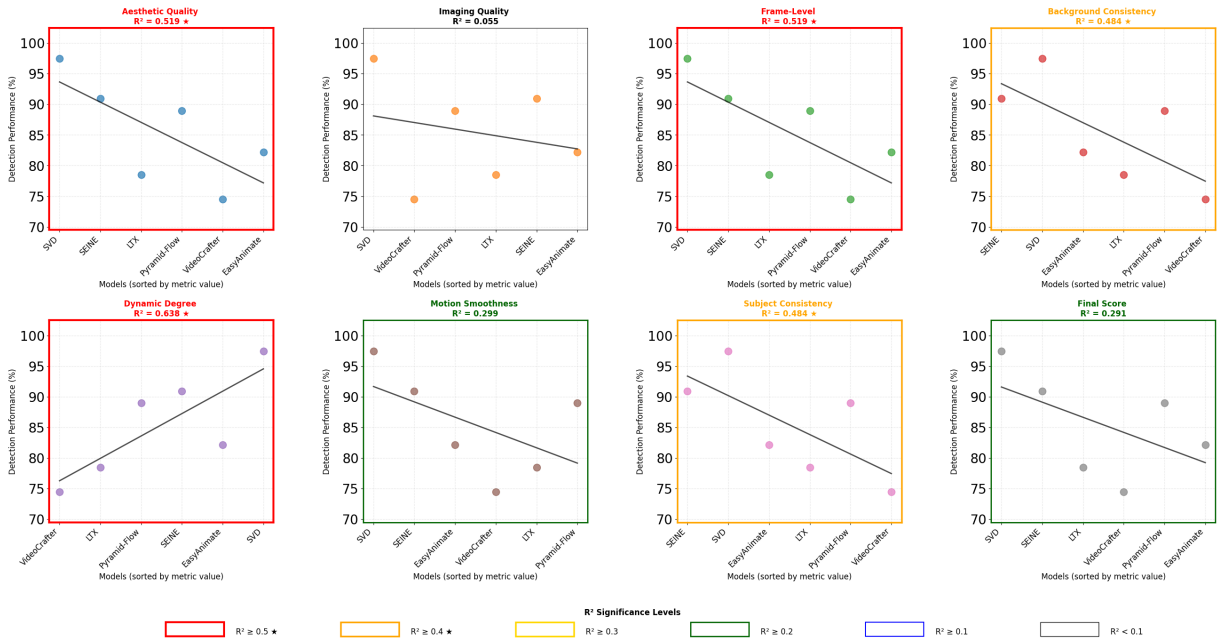


(b) Image-to-Video (I2V) Task.

Figure 14. Correlations between the performance of video classification models and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

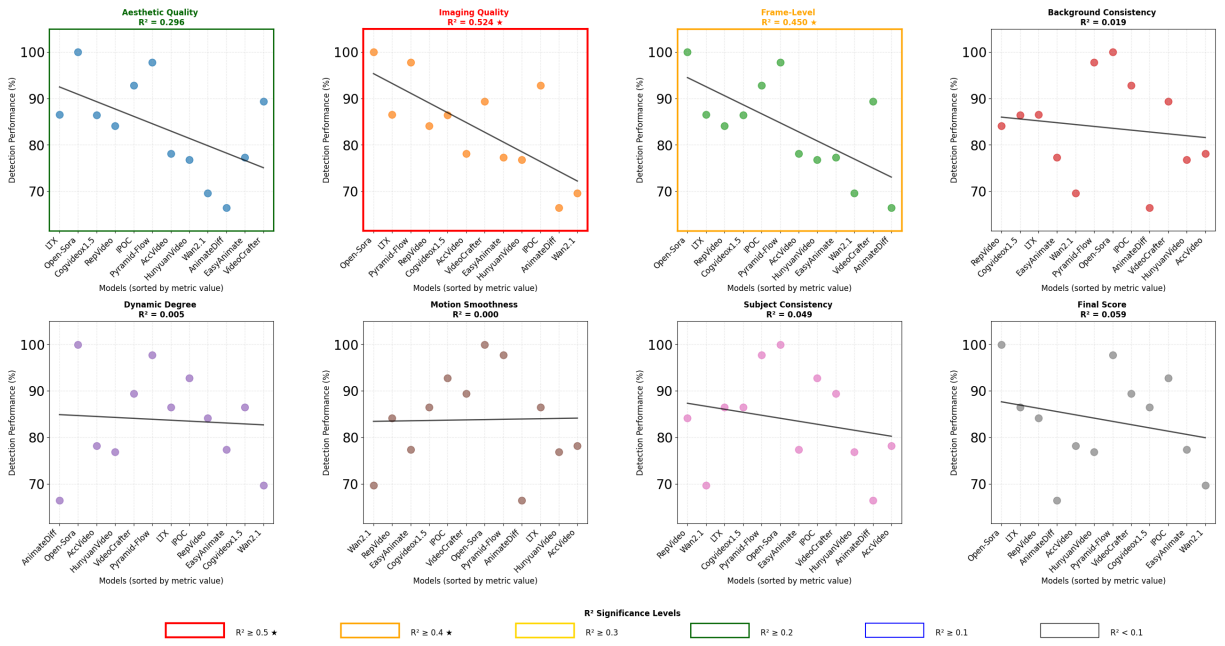


(a) Text-to-Video (T2V) Task.

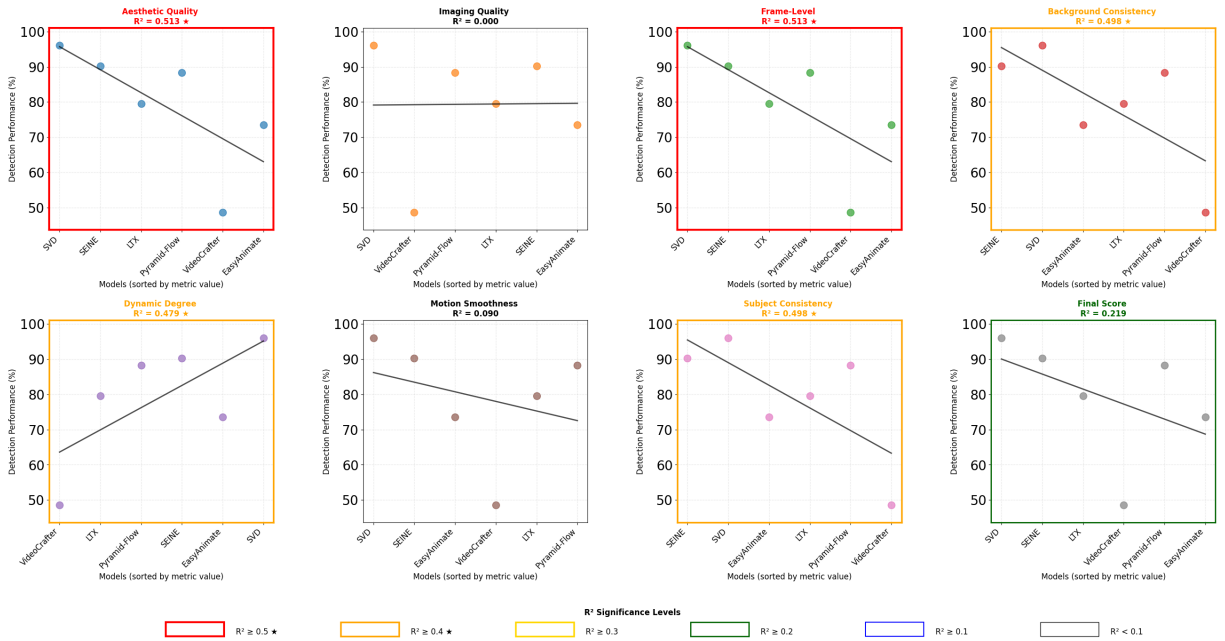


(b) Image-to-Video (I2V) Task.

Figure 15. Correlations between the performance of AI-generated image detection Models and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

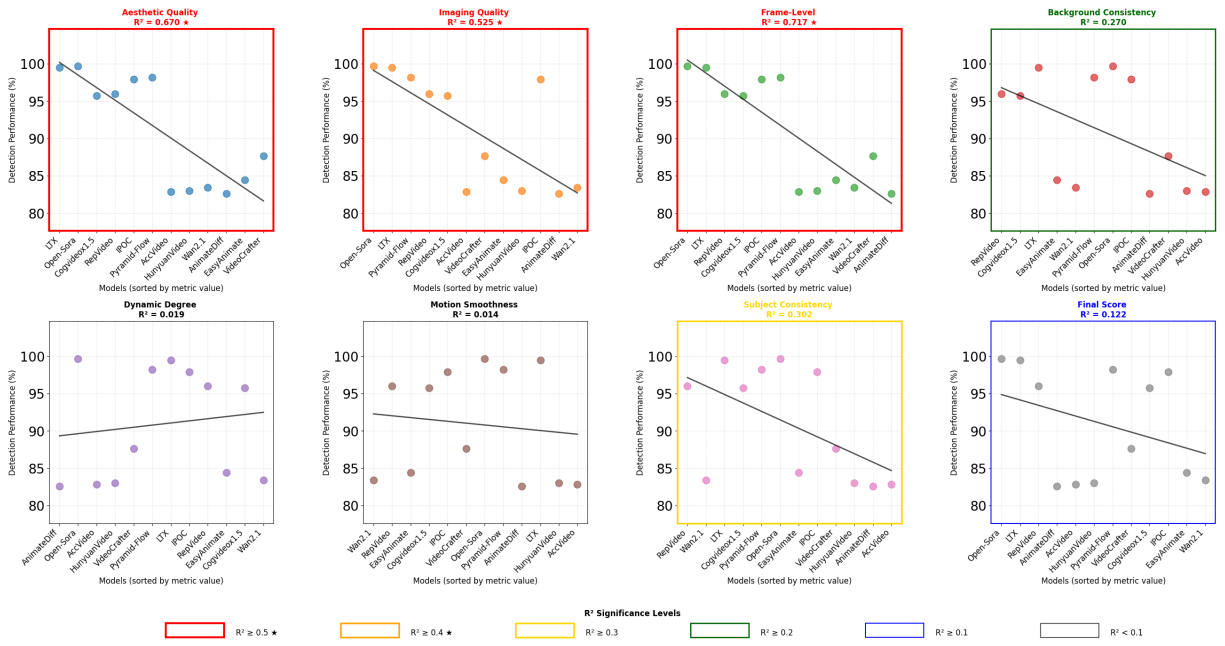


(a) Text-to-Video (T2V) Task.

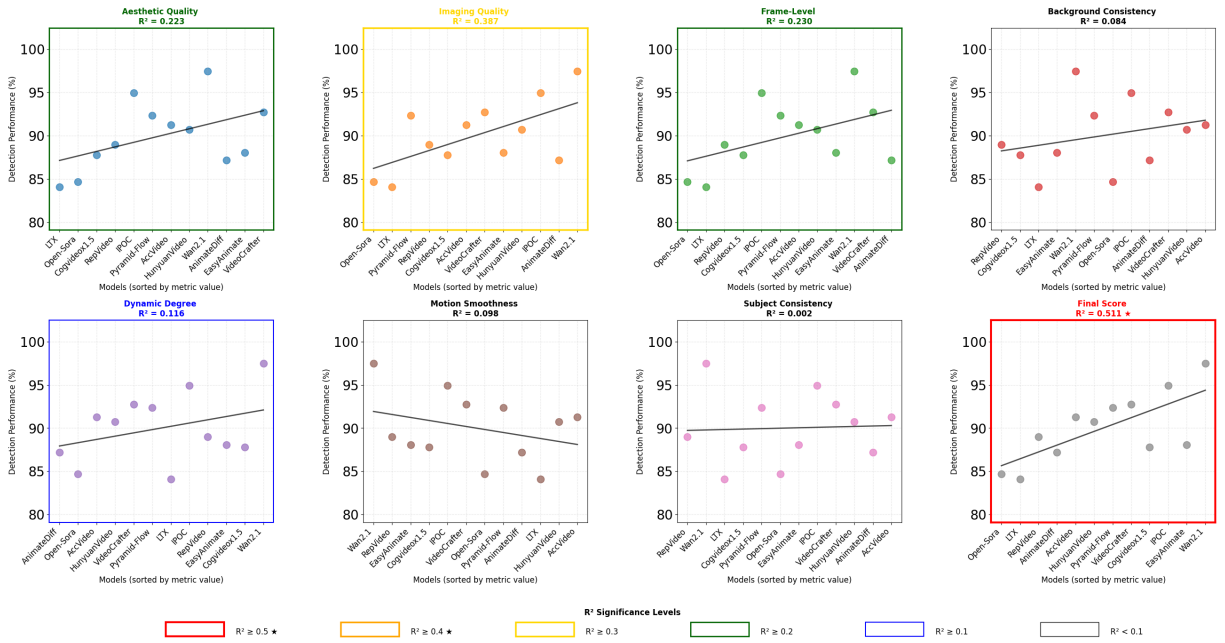


(b) Image-to-Video (I2V) Task.

Figure 16. Correlations between the performance of AI-generated video detection models and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

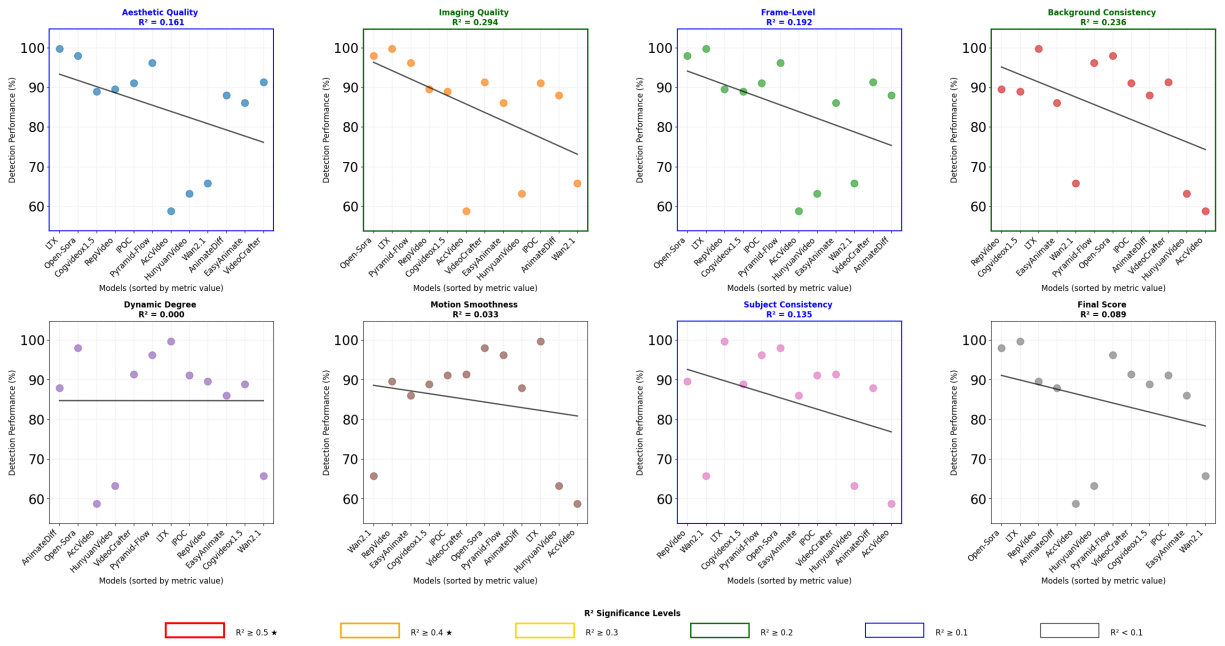


(a) Trained on LTX.

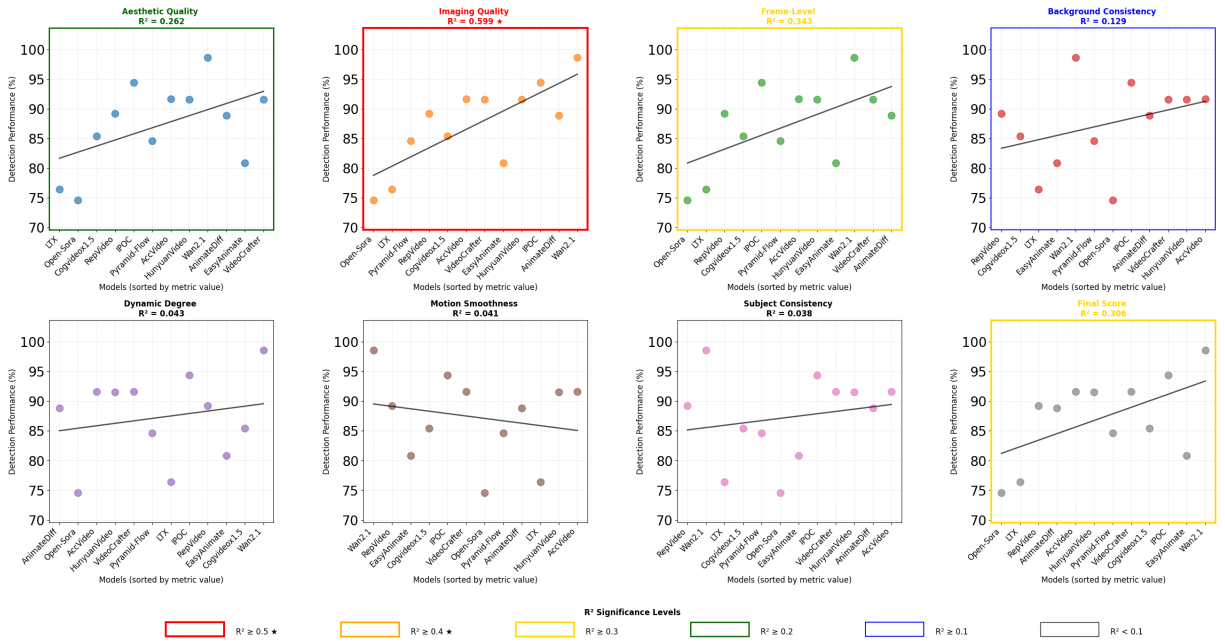


(b) Trained on Wan2.1.

Figure 17. Correlations between the performance of DeCoF and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

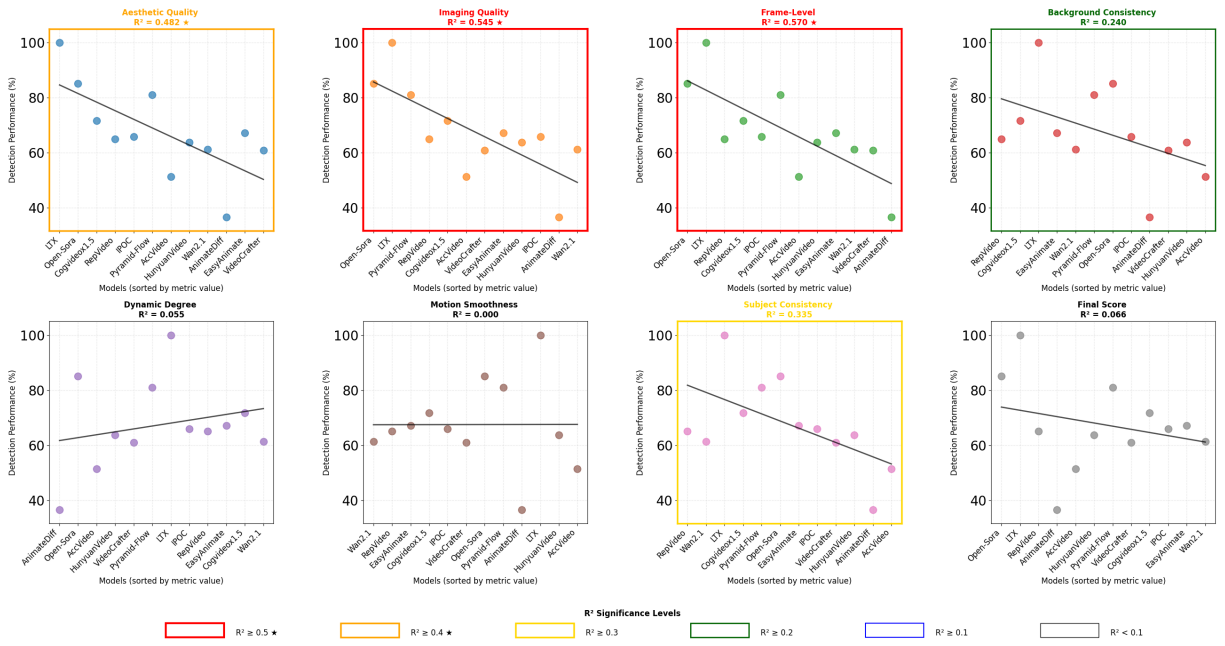


(a) Trained on LTX.

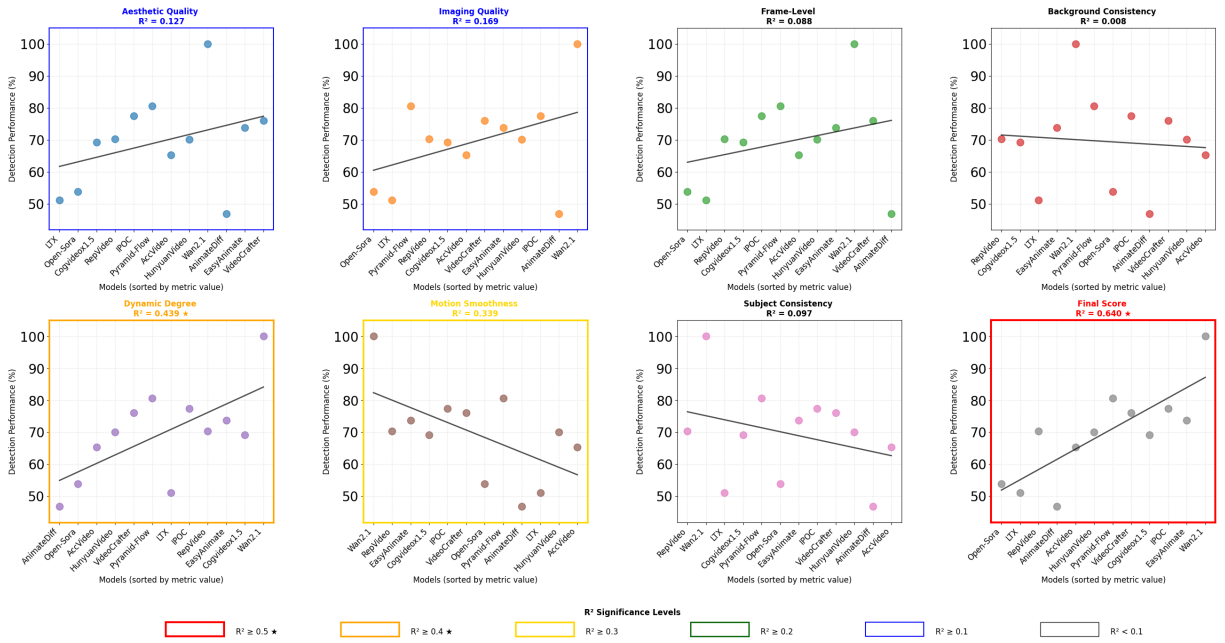


(b) Trained on Wan2.1.

Figure 18. Correlations between the performance of UnivFD and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

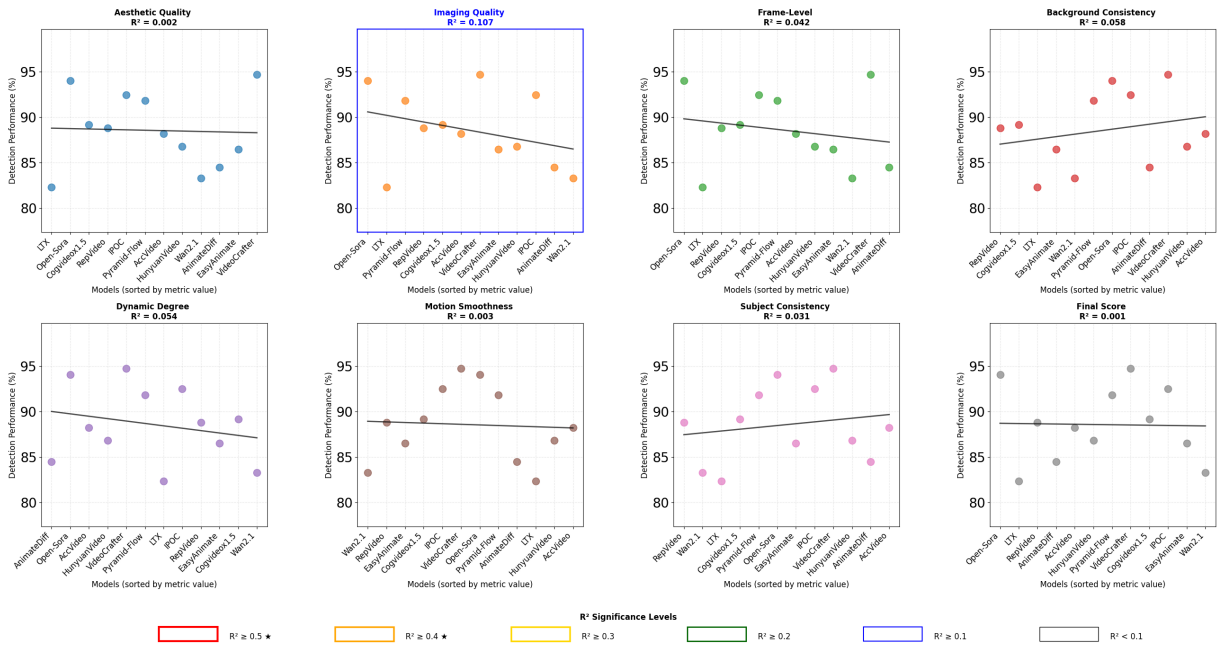


(a) Trained on LTX.

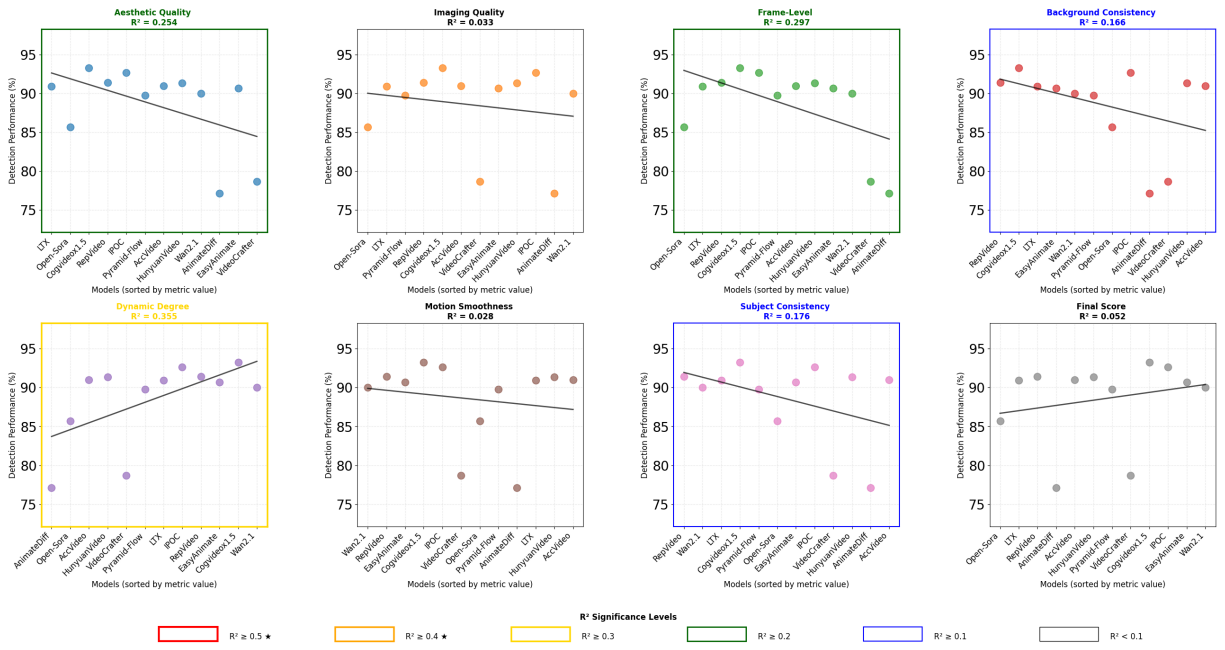


(b) Trained on Wan2.1.

Figure 19. Correlations between the performance of I3D and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

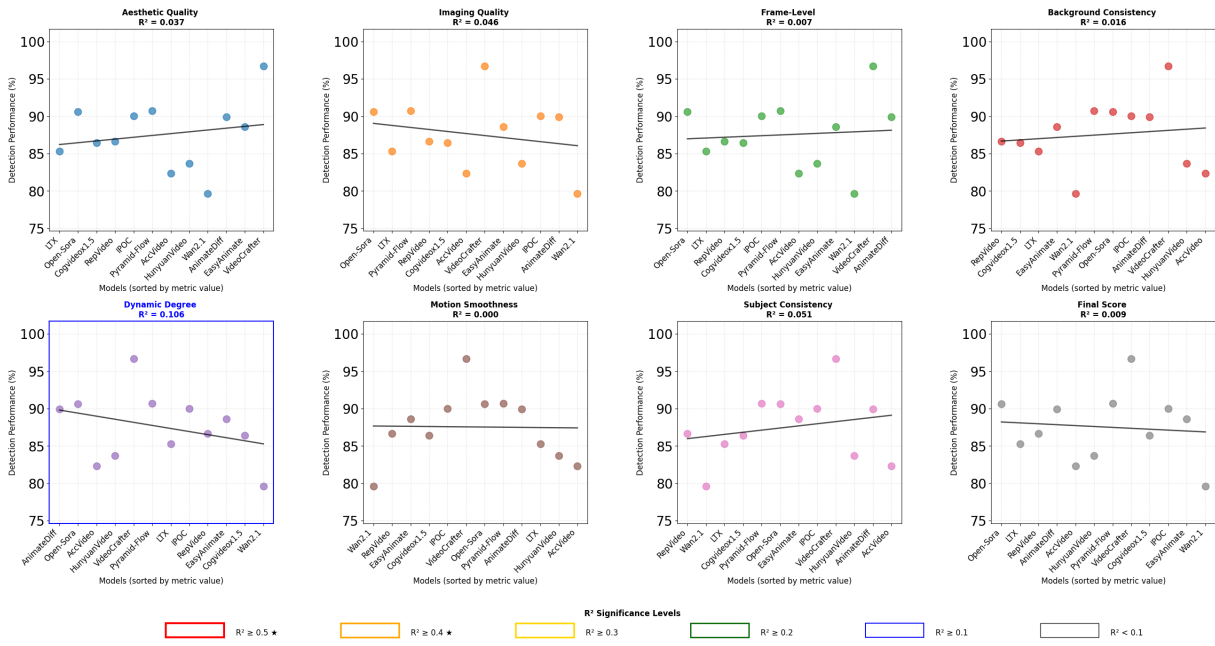


(a) Higher Generative Quality Does Not Increase Detection Difficulty.

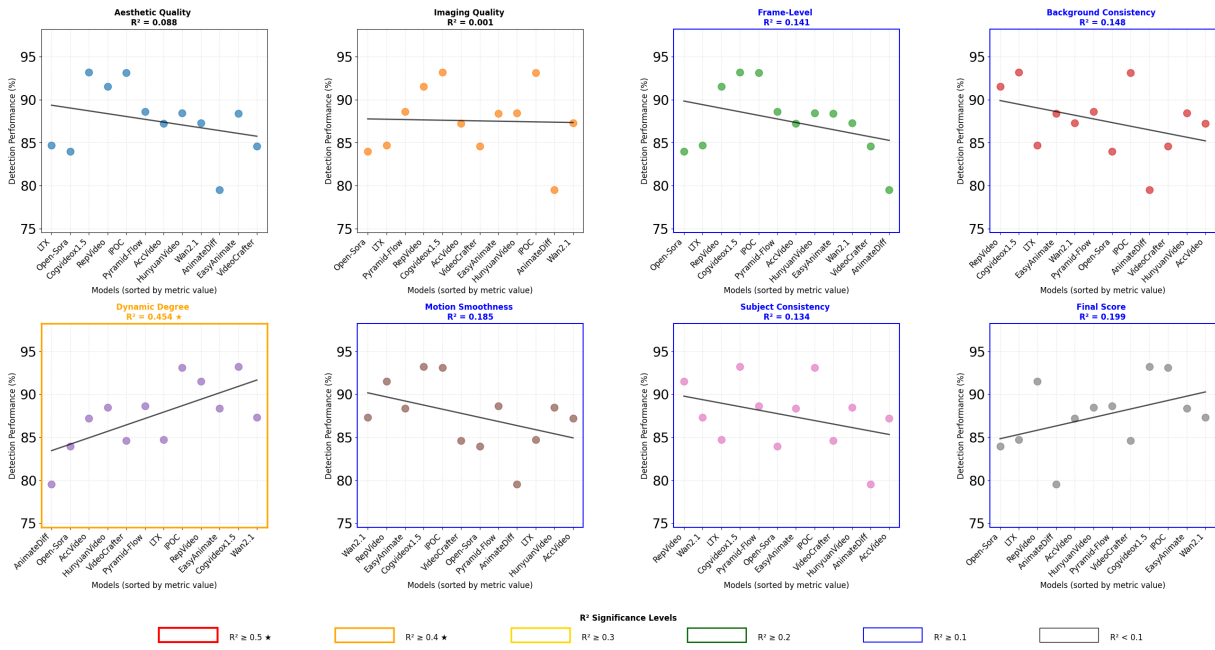


(b) Better Generative Models Do Not Yield Better Detectors.

Figure 20. Validation Through DeCoF: On the Relationship between Generative Model Quality and Detection Efficacy.

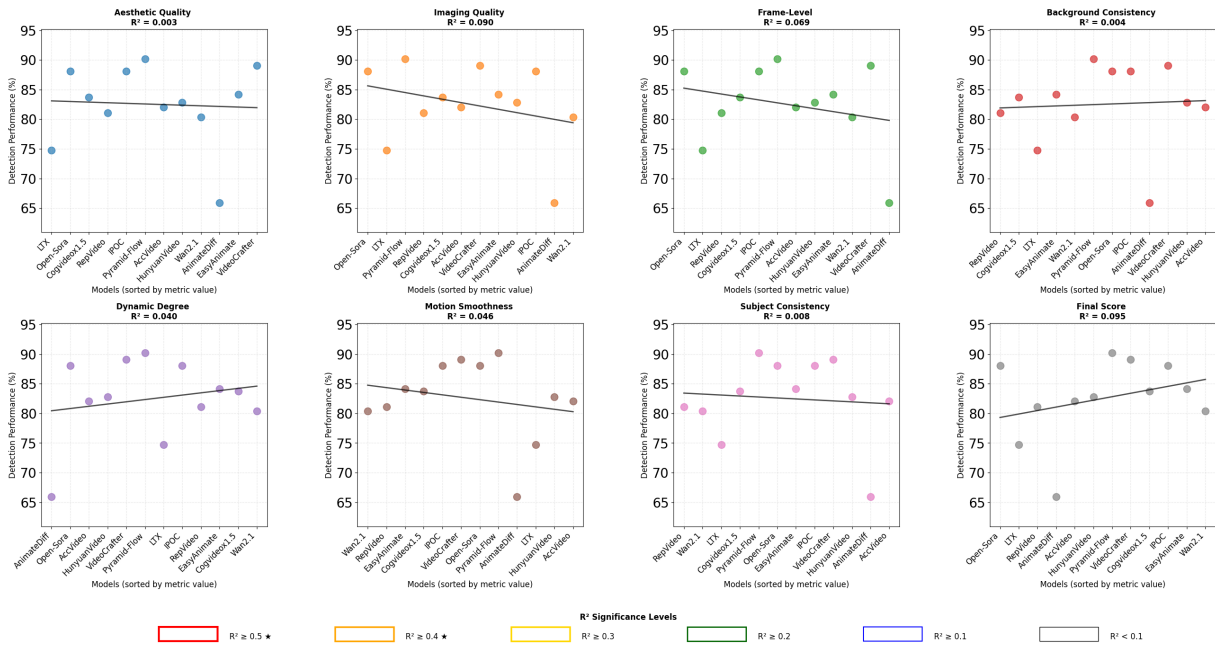


(a) Higher Generative Quality Does Not Increase Detection Difficulty.

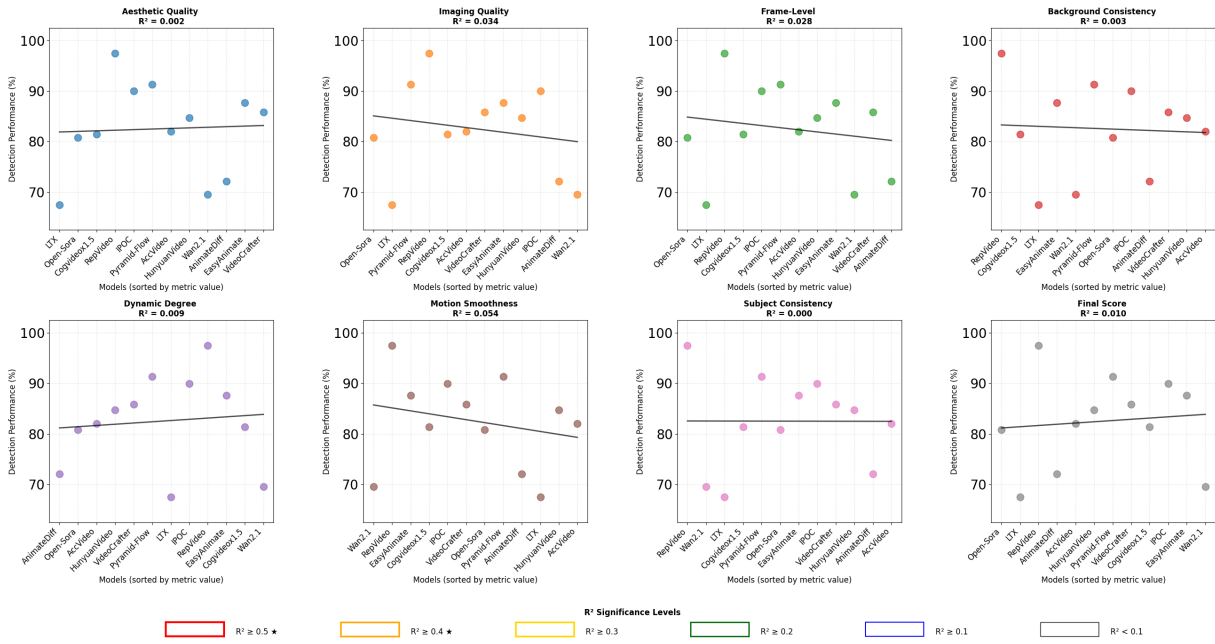


(b) Better Generative Models Do Not Yield Better Detectors.

Figure 21. Validation Through UnivFD: On the Relationship between Generative Model Quality and Detection Efficacy.

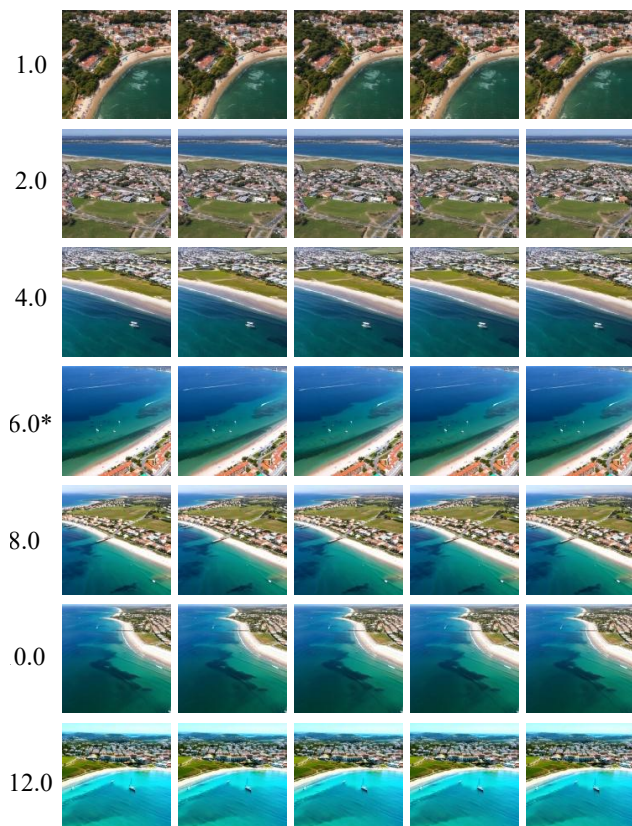


(a) Higher Generative Quality Does Not Increase Detection Difficulty.

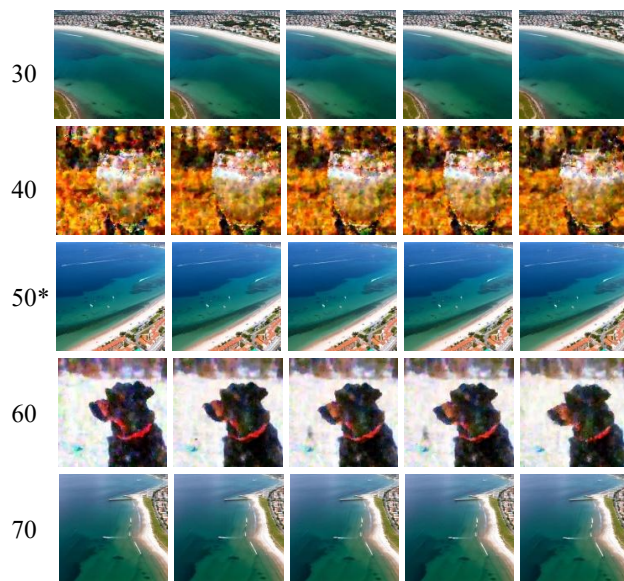


(b) Better Generative Models Do Not Yield Better Detectors.

Figure 22. Validation Through I3D: On the Relationship between Generative Model Quality and Detection Efficacy.



Prompts: an aerial view of a coastal town, showcasing its beachfront and surrounding landscape. ... The overall impression is one of a peaceful, picturesque coastal town.



Prompts: an aerial view of a coastal town, showcasing its beachfront and surrounding landscape. ... The overall impression is one of a peaceful, picturesque coastal town.

Figure 23. The impact of different sampling steps and guidance scale on videos. * Represents the default parameters of the model.

tency across generated videos. This increased uniformity may consequently reduce the challenge of distinguishing between real and synthetic content, thereby improving detection performance for most models.

Analysis-2.5: The impact of video type on detection difficulty.

We have consistently emphasized the importance of dataset scale and diversity. While our previous analysis has demonstrated the significance of scale, we further investigate how different types of videos generated by the same model affect detection difficulty. All videos in this experiment are sourced from the AIGVDBench test set. Since a single video may possess multiple attributes, making clear categorization challenging, we utilize videos with four distinct spatial content types to illustrate this issue. Each category contains 400 videos (200 real and 200 generated), with the spatial content strictly limited to one of the four types. It is also important to note that these videos might still contain spatial elements beyond the four specified categories, as video content is inherently complex, particularly given that we employed detailed textual prompts to control the content generation.

Furthermore, we selected four representative network architectures for evaluation. For instance, CNNSpot employs a ResNet-50 backbone, a architecture widely adopted in generated image detection models. Another example is UnivFD, which utilizes a CLIP-ViT-L/14 framework. For video classification, we included I3D as a representative architecture. As shown in Tab. 20, for lower-performing generative models, the impact of video content variation on detection performance is relatively minimal. However, as the quality of generative models improves, this performance discrepancy becomes increasingly pronounced. Ultimately, the same detector can exhibit performance variations of up to approximately 15% across different video types, underscoring the critical importance of maintaining balanced content representation in datasets.

8. Details of AIGVDBench

8.1. Brief Introduction of Generation Models

A. Video Generation Models. The development of video generation models has undergone a significant evolution from basic generative capabilities to the production of high-quality, long-sequence videos. Early research [8, 24, 33, 49, 95, 100] primarily focused on generating short video clips (around 2–3 seconds), emphasizing fundamental inter-frame consistency and adherence to simple textual prompts. This was later extended to multimodal tasks [14, 20, 60, 63], though the emphasis remained largely on surface-level realism. In recent years, models like Sora [57, 58], kling [35], Gen-3 [67], and HunyuanVideo [34] have shifted focus toward intrinsic realism. Leveraging large-scale training

and Transformer-based architectures, these models achieve long-sequence generation with attention to physical laws, commonsense reasoning, and creative composition, driving video generation toward applications such as AI-assisted filmmaking and world model simulations. However, as AI-generated videos proliferate across social platforms, society is facing a new and far-reaching trust crisis.

B. Used Video Generation Models in AIGVDBench .

During the construction of AIGVDBench, we faced the objective constraint of limited GPU memory resources. Accordingly, the generative models were selected based on the following criteria: the model must be able to run within 80GB of GPU memory while generating a single video in under two minutes, or alternatively, operate within 24GB of GPU memory with a generation time of less than eight minutes per video. Despite these constraints, the dataset construction process still required nearly four months to complete; the full timeline is provided in Sec. 8.8. In the following, we briefly introduce the video generation models used in AIGVDBench, omitting some models that were applied repeatedly.

• Closed source video generation model.

- **CausVid** (CVPR, 2025) is a fast causal video generation model jointly developed by MIT and Adobe. It employs an innovative causal transformer architecture and knowledge distillation to generate high-quality videos with very low initial latency (1.3s for the first frame) and supports streaming generation with real-time instruction modification, excelling in speed and temporal coherence. *At the time of AIGVDBench’s completion, this model was not yet open-sourced and was therefore categorized as proprietary.*
- **Gen2** is a diffusion-transformer video model that turns text prompts or stills into short clips, lets users inject camera trajectories plus segmentation/depth controls, and wraps everything in a creator-friendly web workflow with asset libraries and multi-track timelines for agile ad, music, or narrative projects.
- **Gen3** builds on the Gen2 stack with tighter text-visual alignment, large-scale live-action data, and physics-aware finetuning to deliver longer, steadier shots; it adds motion-reference uploads, mocap presets, and collaborative workspaces aimed at cinematic grading, dialogue-driven scenes, and bespoke brand content.
- **Kling** is Kwai’s multimodal video suite built on a 3D diffusion transformer with spatio-temporal latents, covering text-to-video up to roughly a minute, image-to-video expansion, character animation, physics priors, control channels, and audio alignment so creators can preserve identity consistency, steer camera choreography, and rapidly deliver music videos, ad spots, or concept reels.
- **Luma Dream Machine** is Luma Labs’ generative video and 3D creation suite, blending diffusion backbones with neural radiance field models to turn text, images, or captured footage into high-res videos, panoramas, and editable NeRF assets, and shipping browser-based tools for lighting, material tweaks, and camera path planning so teams can iterate on commercials, concept art, or digital twins with minimal friction.
- **Open-Sora** is an open-source project initiative aimed at replicating and advancing the capabilities of OpenAI’s Sora model. It focuses on significantly reducing the training costs associated with high-quality video generation models and promoting the development of accessible, open-source alternatives in the field, with its v2.0 version showing performance close to the closed-source Sora model.
- **Pika** is an AI video generation tool developed by Pika Labs, enabling users to create and edit videos from text prompts. It is recognized for its user-friendly interface and ease of use, making it popular for social media content creation and quick video prototyping, allowing for efficient creative expression.
- **Sora** is a groundbreaking text-to-video model developed by OpenAI. It is capable of generating high-fidelity, temporally coherent video clips up to one minute long directly from text descriptions, setting a notable benchmark for realism and narrative complexity in AI-generated video and demonstrating significant advancements in the field.
- **Vidu** is an AI video generation model capable of producing videos from text descriptions. Its capabilities have been demonstrated in public showcases, highlighting its potential in the text-to-video generation space and indicating its role as a contender in the rapidly evolving field of AI video synthesis.
- **Wan** is a comprehensive and open video foundation model suite introduced by Alibaba’s Wanxiang team. It includes models with 1.3B and 14B parameters, covering multiple downstream tasks like text-to-video and image-to-video generation. Notably, it is the first model capable of generating visual text in both Chinese and English.

• Open source video generation model.

- **EasyAnimate**: A high-performance, open-source long video generation project developed by Alibaba Cloud’s PAI team, based on a Diffusion Transformer (DiT) architecture. It modifies the VAE and DiT structures to better support video generation, extending the original 2D image synthesis DiT framework into the 3D video domain and introducing a motion module to capture temporal dynamics. The proposed Slice VAE addresses GPU memory challenges during the encoding and de-

Table 21. AIGVDBench Used Video Generation Models Configuration.

Model	TeaCache	Resolution	Num Frames	Num Steps	CPU Offload	Precision	Guidance Scale	FPS
Open-Sora	None	256 (16:9)	37	50	True	bf16	7.5	24
Hunyuan	0.15	544 × 960	33	50	True	fp8	6.0	24
Wan2.1	0.08	480 × 832	49	50	True	bf16	5.0	16
AccVideo	None	544 × 960	33	5	True	fp16	1.0	24
AnimateDiff-L	None	512 × 512	16	8	False	—	1.0	10
CogVideoX	None	480 × 720	81	50	False	fp16	6.0	16
Pyramid-Flow	None	384 × 640	89	10+20	False	bf16	7.0	24
VideoCrafter	None	320 × 512	16	50	False	—	12.0	10
EasyAnimate	0.08	384 × 672	49	50	True	bf16	6.0	8
LTX(T2V)	0.05	512 × 768	101	50	False	bf16	3.0	24
LTX(I2V)	0.05	513 × 768	50	50	False	bf16	3.0	24
LTX(V2V)	None	514 × 768	50	50	False	bf16	3.0	24
RepVideo	None	480 × 720	49	50	False	fp16	6.0	8
SEINE	None	320 × 512	16	250	False	fp16	8.0	8
SVD	None	320 × 513	16	25	False	—	1.0-3.0	5
IPOC	None	480 × 720	49	50	False	bf16	6.0	8

- **TeaCache:** TeaCache [43] acceleration (enabled/disabled) and the associated *rel_11_thresh* parameter.
- **Offload:** Whether the model employs memory offloading techniques.
- **FPS:** Frames per second in the generated video output.
- **10+20:** *num_inference_steps* : [20, 20, 20], *video_num_inference_steps* = [10, 10, 10].

coding of long, large videos, enhancing compression efficiency along the time dimension. It offers a complete solution for HD long video generation, including data preprocessing, VAE training, DiT training, model inference, and evaluation.

- **LTX:** An all-in-one generative AI platform designed for video production, catering to filmmakers, advertisers, and creative teams. It transforms creative ideas into high-quality, professional videos, streamlining stages from scripting and storyboarding to editing and final delivery. Its AI creative engine, LTX-2, powers synchronized audio and video generation, offers 4K fidelity, and includes multiple performance modes for efficiency. The platform provides tools for controlling camera movement, defining motion with keyframes, using visual references, and maintaining consistency with AI characters / objects.
- **Pyramid-Flow** (ICLR 2025): An open-source high-definition video generation model jointly released by Peking University, Beijing University of Posts and Telecommunications, and Kuaishou Technology. It employs a pyramid flow matching algorithm, which decomposes the video generation process into multi-

ple stages starting from low resolution and progressively upgrading to high resolution. This approach effectively handles high-dimensional video data, significantly reducing computational costs and token count while maintaining high visual quality.

- **SEINE** (ICLR 2024): A diffusion model designed for generative transitions and prediction, focusing on generating videos from short to long durations. As part of the Vchitect video generation system, it can be integrated with text-to-video frameworks like LaVie. Based on Stable Diffusion v1.4, SEINE is particularly adept at creating smooth transitional effects and predicting video content, making it suitable for applications in filmmaking, advertising, and video editing.
- **SVD:** A pioneering video generation model that applies the diffusion transformer architecture to video synthesis, developed by Google Research. It is capable of generating short video clips with coherent motion directly from a single static image, effectively animating still pictures. The model demonstrates significant capabilities in producing realistic and temporally consistent motions from image inputs, establishing a strong baseline for image-to-video generation research.

- **VideoCrafter** (CVPR 2024): An open-source video generation model and toolkit, with versions like VideoCrafter2 supporting high-quality text-to-video and image-to-video generation. It aims to provide a reproducible and modifiable solution for the video generation community, often serving as a strong baseline for research and application development. The project typically offers a complete ecosystem, including model training, inference, and evaluation tools.
- **AnimateDiff**: An influential open-source framework designed to inject motion dynamics into personalized text-to-image models, enabling them to generate video clips. Its core innovation lies in learning a universal motion module that can be plugged into various pre-existing personalized T2I models without modifying their original weights. This approach efficiently animates a wide range of personalized concepts, significantly advancing personalized video generation.
- **CogVideoX** (ICLR 2025): A powerful open-source video generation model developed by Tsinghua University. Based on the Transformer architecture, it supports text-to-video and image-to-video generation, as well as video continuation. The model employs technologies like 3D VAE and 3D RoPE positional encoding to significantly enhance temporal coherence and long-range dependencies in the generated videos.
- **HunyuanVideo**: A sophisticated text-to-video generation model developed by Tencent. It leverages a robust multi-stage training strategy, progressing from low to high resolution and from image pre-training to video training, to achieve high-fidelity output. A key strength is its deep integration of the proprietary Hunyuan large language model for enhanced text comprehension and semantic alignment.
- **IPOC**: A research direction in video generation focused on complex scene composition and open-world synthesis from image prompts. It tackles challenges like seamlessly integrating multiple elements from reference images into a cohesive and dynamic video scene, aiming for greater creative control and visual consistency.
- **RepVideo**: A video generation method that emphasizes efficient and high-fidelity reproduction of visual elements. The approach focuses on achieving high-quality replication of details, motions, or styles within the generated video content through innovative architectural designs or training strategies.
- **Wan2.1**: A version iteration of the "Wan" video generation model suite, potentially from Alibaba's Wanxiang team. This version introduces improvements over its predecessors in areas such as generation quality, supported resolutions, and specific capabilities like

generating visual text in multiple languages. It is designed to be an efficient model suitable for consumer-grade hardware.

- **AccVideo**: A novel extraction method focused on accelerating existing video diffusion models. Its core innovation lies in using a pre-trained model to generate a high-quality synthetic dataset. Through trajectory-based few-step guidance and an adversarial training strategy, the student model learns to produce high-quality videos with very few inference steps, achieving significant speed boosts while maintaining comparable video quality.

The detailed parameters of the generative models can be found in Tab. 21.

8.2. Introduction of the Used Detection Methods

A. Generated Image Detection. Existing approaches for synthetic image detection can be broadly categorized into two main research directions. The first line of methods focuses on forgery artifact learning, which aims to capture low-level traces left by generative models in either the spatial or frequency domain. For instance, some studies employ backbone networks such as ResNet combined with specialized modules to extract frequency features [23], gradient information [45, 70], or reconstruction anomalies [80]. Other works analyze local pixel relationships arising from upsampling operations [71]. Although these methods perform well on known generative models, their generalization is often constrained by the training data distribution, and they typically require large-scale training from scratch. The second category builds on frozen pre-trained vision models, such as CLIP-ViT-L/14, to extract general-purpose features, with only a lightweight classification head being fine-tuned for detection [56]. Such methods mitigate overfitting to specific artifacts and demonstrate promising cross-model generalization in few-shot scenarios. However, the frozen encoders often retain a substantial amount of forgery-irrelevant information, limiting further performance gains. Recent efforts have sought to incorporate forgery-aware priors by introducing prompt tuning or adapter modules into frozen architectures [17, 88], advancing this line towards more efficient and discriminative detection.

B. Deepfake Detection. The field of deepfake face detection has evolved to address the unique challenges of media forgery, with research directions broadly classified into the following categories: 1. Detection Based on Inherent Forgery Features: Early methods primarily leveraged convolutional neural networks (CNN) as end-to-end binary classifiers [1, 65]. Subsequent work has specialized by targeting specific forgery traces. Spatial-domain methods focus on cues such as blending boundaries [38, 97] or localized artifacts [54]. 2. Detection Enhanced by Data Aug-

mentation: To improve generalization, a prominent strategy involves augmenting training data with synthetic forgeries. This includes generating pseudo deepfakes by applying face warping [38], self-blending [13, 68], or adversarial blending simulations [87]. The core idea is to preemptively expose the detector to a wider variety of forgery types and artifacts, thereby enhancing its robustness against unseen manipulation techniques. 3. Detection Focused on Generalization Capabilities: Given the rapid evolution of generation models, a significant research thrust is dedicated to improving cross-dataset and cross-method generalization. Efforts span multiple dimensions, including advanced data augmentation [51, 85, 86], frequency feature analysis [19, 64], and network architecture innovations [97].

C. Used Detection Methods in AIGVDBench. We categorize potential approaches for synthetic video detection into four groups: (1) video classification models, which employ specialized network architectures for video based tasks; (2) AI generated image detection models, which can be adapted as frame level detectors for synthetic videos; (3) dedicated AI generated video detection models specifically designed for this task; and (4) the transfer of vision language models, which also represents a highly promising direction. Based on this taxonomy, we selected 33 mainstream models for evaluation. It should be noted that Deepfake detection models were excluded from our study, as they are designed specifically for facial forgery and are not suitable for general synthetic video detection. Additionally, methods requiring extensive input preprocessing, such as video reconstruction, were also omitted due to the substantial scale of our dataset, where such operations would significantly impede evaluation efficiency. In the following section, we introduce the 33 detectors included in our benchmark.

• **Video Classification models.**

- **MViT (CVPR, 2022):** An open-source repository from Facebook Research (Meta AI) that implements Multiscale Vision Transformers (MVIT), a model designed to address limitations in standard Vision Transformers regarding computation and multiscale feature modeling. MViT introduces a hierarchical transformer architecture that combines self-attention with pooling operations, enabling efficient capture of features at multiple scales. This leads to improved performance and reduced computational cost compared to conventional models. The repository provides modular PyTorch implementations for image classification, object detection, and video understanding tasks, achieving state of the art results on benchmarks such as ImageNet [18] and Kinetics [32].
- **UniFormer (ICLR, 2022):** A unified transformer framework that seamlessly integrates convolution and self-attention into a single network architecture for both video understanding and image recognition. It delivers state of the art performance by merging local and global modeling strengths, and features flexible scalability for various vision tasks.
- **VideoSwin (CVPR, 2022):** A video understanding model based on the Swin Transformer architecture, which efficiently captures spatio-temporal features in videos through a hierarchical window based attention mechanism. It is widely used in tasks such as action recognition and video classification, demonstrating excellent performance.
- **VideoMAE (NeurIPS, 2022):** A self-supervised video pre-training model that employs Masked Autoencoders (MAE) to efficiently learn spatiotemporal features from videos. It randomly masks a large portion of space-time blocks in the input video and reconstructs the full video from the remaining visible parts. This encourages the model to learn richer and more robust representations, drastically improving data efficiency. VideoMAE achieves state of the art performance in various video understanding tasks (like action recognition), advancing self-supervised learning in the video domain.
- **TSM (ICCV, 2019):** An efficient video understanding model that introduces the "Temporal Shift Module" for effective temporal modeling. This approach shifts a portion of feature channels along the temporal dimension, enabling information exchange between frames without introducing additional parameters or computation. TSM achieves a good balance between efficiency and performance, making it suitable for various video understanding tasks like action recognition, and it has achieved excellent results on several benchmark datasets.
- **SlowFast (ICCV, 2019):** utilizes dual pathways: a "slow" pathway and a "fast" pathway, capturing video spatiotemporal features at different temporal resolutions for precise understanding of dynamic visual content. The codebase supports tasks such as action recognition, making it widely used in both research and practical applications.
- **TimeSformer:** is an end to end video understanding model based purely on Transformers, applying attention mechanisms along both spatial and temporal dimensions to effectively extract video features. The model has driven the adoption of pure self-attention methods in tasks like video classification.
- **I3D (CVPR, 2017):** A convolutional neural network model for video classification based on the I3D (Inflated 3D ConvNet) architecture, trained on the large scale Kinetics human action recognition dataset. I3D

extends 2D convolutions to 3D, effectively modeling spatiotemporal features for better recognition of dynamic actions in videos. This project provides a strong baseline for video understanding tasks and is widely used in action recognition research.

- **UniFormerV2 (ICCV, 2023)**: A state of the art model for spatiotemporal learning in videos. UniFormerV2 innovatively combines the Image Vision Transformer (ViT) architecture with the video specific UniFormer framework, leveraging the strengths of both for spatial and temporal modeling. This model efficiently captures rich spatiotemporal information within videos, enhancing performance in tasks like action recognition, and offers a more unified and general solution for video representation learning.
- **X3D (CVPR, 2020)**: A deep learning model for efficient video recognition. X3D expands existing 3D convolutional neural network architectures along multiple axes—such as spatial, temporal, and network width—achieving excellent video understanding performance with minimal computational cost. The model is lightweight and fast, making it well suited for action recognition and other tasks in resource constrained settings, and it has demonstrated strong results on various video benchmarks.
- **AI-Generated Image Detection Models.**
 - **Fredect (ICML, 2020)**: is designed to automatically identify AI-generated or manipulated content in images or videos. It integrates multiple feature analysis methods and deep learning techniques to detect anomalous patterns or artifacts in visual data, thereby improving the recognition accuracy of synthetic media. Fredect is applicable in academic research, content moderation, and digital forensics, serving as a key technology for enhancing media authenticity and security.
 - **NPR (CVPR, 2024)**: A deepfake detection project leveraging Non-Photorealistic Rendering (NPR) features. By extracting NPR-based image representations, this method effectively discriminates real from fake content, improving robustness and generalization across diverse generators and complex backgrounds. The project is well-suited for applications in digital media security and content authentication.
 - **Fusing**: A project focused on fusing global and local features for video understanding tasks. It aims to jointly model high-level contextual information and fine-grained local details in videos, enhancing the representation of complex dynamic content and improving video analysis accuracy. Such fusion strategies are valuable in domains such as action recognition and multimodal reasoning, supporting more precise interpretation of video content.
- **Gram-Net (CVPR, 2020)**: By enhancing global texture features in facial images, the model improves the accuracy and robustness of forged face detection, particularly against DeepFake and face-swapping techniques. Combining texture analysis with deep learning, the method advances fake face detection under real-world conditions and is widely applicable in media authentication and cybersecurity.
- **CNNspot (CVPR, 2020)**: It detects CNN-synthesized images by identifying subtle statistical and textural artifacts specific to synthetic images. The method effectively distinguishes real from AI-generated content, demonstrating strong performance on early generative models and providing valuable insights for deepfake detection and image forensics.
- **D3 (CVPR, 2025)**: A Discrepancy Deepfake Detector designed for multi-generator deepfake detection. It employs a dual-branch architecture: a main branch that processes original image features and an auxiliary branch that handles specially distorted inputs. By fusing features from both branches, the model captures common discrepancy signals across generators and decomposes complex forgery patterns, leading to improved adaptability and generalization in detecting fake content from diverse sources.
- **ForgeLens (ICCV, 2025)**: It introduces novel network structures and training strategies that effectively localize forged regions within images, achieving strong generalization even with limited annotated data. The method is applicable to deepfake and image manipulation detection, offering an efficient and universal solution for digital media security. Based on distinct experimental configurations, we categorized the evaluation settings into ForgeLens1 and ForgeLens3.
- **Effort-AIGI-Detection (ICML, 2025)**: This work improves the generalization of AI-generated image detection via Orthogonal Subspace Decomposition (OSD), which decomposes image features into independent components to better separate real from fake characteristics. The approach significantly enhances cross-domain and multi-generator detection performance, making it suitable for real-world applications in digital media security and deepfake forensics.
- **AI-Generated Video Detection Models.**
 - **DeMamba** focus on large-scale AI-generated video detection. In response to the rapid growth of video generation technology and the popularity of video content on social media, this project aims to develop methods and models capable of distinguishing between real and AI-generated videos over million-scale datasets (Gen-Video Benchmark). The DeMamba framework extracts and analyzes features from both genuine and fake

videos, helping curb the spread of misinformation and meeting the rising demand for efficient and accurate video forgery detection tools, especially for digital media security and content moderation use cases.

- **DeCoF** focus on detecting AI-generated videos via frame consistency analysis. By examining the consistency between video frames and leveraging temporal and sequential features, DeCoF effectively identifies AI-generated content and uncovers subtle artifacts that may not be apparent in individual frames. This approach enhances detection performance for complex generative videos and is suitable for applications such as media content moderation and video forensics, helping to limit the spread of fake video information.
- **Vision-Language Models.**
 - **DeepSeek-VL-7B-Chat:** An instruction-tuned multimodal model that pairs SigLIP-L and SAM-B encoders for 1024×1024 visual input with the DeepSeek-LLM-7B backbone. It was trained on approximately 2 trillion text tokens and 400 billion multimodal pairs. The model features conversational formatting, strong diagram and document comprehension capabilities, and flexible image grounding APIs, making it suitable for general-purpose assistant applications.
 - **InternVL3-8B:** Integrates an InternViT-300M vision encoder with a Qwen2.5-7B decoder through native multimodal pretraining, mixed preference optimization, and variable visual position encoding. The model targets diverse applications including document analysis, GUI understanding, spatial reasoning, video processing, and multilingual tasks, while providing tool-calling capabilities and video-aware chat templates.
 - **Qwen2.5-VL-3B-Instruct:** A lightweight model combining a Vision Transformer front end (featuring windowed attention, SwiGLU, and RMSNorm) with a 3B Qwen2.5 decoder via a multimodal projector. Dynamic spatial and temporal sampling enables unified image and video processing, while the compact decoder ensures fast inference and precise instruction following.
 - **Qwen2.5-VL-7B-Instruct:** Scales the ViT-projector-LLM architecture to a 7B language core and higher-capacity vision encoder, maintaining dynamic resolution and mRoPE temporal encoding. Structured-output heads provide reliable performance for chart parsing, UI reasoning, and enterprise document extraction pipelines.
 - **Qwen2.5-VL-32B-Instruct:** Expands the modular ViT + projector architecture with a 32B parameter decoder, incorporating reinforcement-tuned preference heads and extended-position embeddings. Capable of

handling 32K+ context windows, multi-image analytics, and complex agent-style planning with enhanced reasoning depth.

- **DeepSeek-VL2:** A Mixture of Experts pipeline combining a SigLIP-style encoder and dynamic tiling with the DeepSeekMoE-27B language expert pool, activating only a subset of experts per token. This design balances latency with strong performance on chart understanding, OCR, and visual grounding, while supporting multi-image conversational prompts.
- **DeepSeek-VL2-Small:** A compact version of VL2 based on the DeepSeekMoE-16B backbone with approximately 2.8B activated parameters. Retains the dynamic tiling preprocessor for up to two images, suitable for memory-constrained deployments while maintaining instruction stability and context-aware visual question answering.
- **Emu3-Stage1:** First-stage Emu3 weights trained exclusively with next-token prediction over tokenized text and images using a single decoder-only transformer with 5120 context length. Leveraging a learned vision tokenizer and classifier-free guidance, it provides unified generation and perception capabilities, serving as the foundation for subsequent video-augmented stages.
- **LLaVA-1.5-7B:** A classic architecture connecting a CLIP ViT-L/14 encoder to a Vicuna-7B decoder through a learned projector. Aligned via MSP/CC-SBU pretraining and GPT-curated VQA instructions, it excels at conversational image descriptions, simple grounding, and few-shot captioning with a well-documented chat template.
- **FastVLM-7B:** Integrates the FastViTHD hybrid encoder—which generates fewer high-value tokens through strided attention—with a Qwen2-7B decoder and slim projector. Engineered for rapid time-to-first-token, it maintains competitive performance on OCR, chart understanding, and document comprehension while providing ready-to-use inference scripts.
- **Kimi-VL-A3B-Instruct:** Combines a MoonViT native-resolution encoder with a 2.8B-activated Moonlight-16B MoE language core via an MLP bridge. Supporting 128K context windows, multi-image and video inputs, and agent-style workflows, it includes specialized prompts for UI automation and long-document analysis.

8.3. Introduction of Original Data

The real video in our AIGVDBench are sourced from the OpenVidHD dataset [53], a large-scale high-quality video collection specifically designed for high-definition video generation research. Derived from OpenVid-1M, it con-

sists of a curated selection of approximately 433,000 high-resolution video clips at 1080p, totaling about 11,000 hours of content. Characterized by stringent quality control, OpenVidHD ensures superior visual quality and diverse thematic coverage, including natural landscapes, human activities, animal behaviors, and more. Each video clip is accompanied by accurate textual descriptions. Designed to address current limitations in high-definition video generation, the dataset provides a valuable foundation for training and evaluating text-to-video generation models, thereby advancing the development of video synthesis technologies. **In the OpenVidHD dataset, all videos are compressed using the MPEG-4 Part 2 format. For consistency in our evaluation, we have standardized the compression to H.264 across all videos in the AIGVBench benchmark.**

8.4. Details of Prompts Categorization

To balance classification accuracy and computational efficiency, we perform prompt categorization based on predefined WordNet synsets and keyword lists. This is achieved by defining a set of WordNet synsets and keywords for each category, then automatically scanning input prompts for words belonging to these lexical resources. When matches are found, the corresponding category labels are assigned to the prompt, enabling multi-label classification.

The WordNet synsets and corresponding keyword lists are presented in Tabs. 22 and 23. It should be acknowledged that, despite our best efforts to ensure classification accuracy, the long-text nature of our prompts inevitably introduced a certain degree of misclassification. Representative examples of such prompts, including incorrectly categorized cases, are provided in Tab. 24. However, given the substantial volume of data (approximately over 400,000 entries) and the considerable length of individual prompts, it was not feasible to employ large language models (LLMs) for this classification task. It is worth noting that we also attempted to utilize closed-source LLMs such as GPT-4, but found that this approach still could not fully resolve the issue.

8.5. Balance Comparison with Prior Datasets

To quantitatively evaluate the balance and diversity of the dataset, we introduced **Penalized Global Balance Score (PGBS)**. At the attribute level, we first introduce the **Relative Uniformity (RU)** to quantify the distribution uniformity of values within a single attribute:

$$RU_i = \frac{H_i}{H_i^{unif}} \quad (1)$$

where denotes the empirical entropy of the i -th attribute, $H_i^{unif} = \log n_i$ represents the entropy under an ideal uniform distribution, and n_i is the total number of theoretical

categories for that attribute. An **RU** value close to 1 indicates that the distribution of the attribute approaches uniformity. Considering that datasets often suffer from missing categories, we extend **RU** to **Complete Uniformity (CU)** as follows:

$$CU_i = RU_i \times (R_i)^\alpha \quad (2)$$

where $R_i = \frac{|S_i^{obs}|}{n_i}$ is the completeness ratio, $|S_i^{obs}|$ denotes the number of actually observed categories, and α is an adjustable penalty intensity parameter. This design ensures that when category missing occurs (i.e., $R_i < 1$), the **CU** value is penalized appropriately, thereby more accurately reflecting the impact of distribution completeness on balance. After obtaining the **CU** values for all attributes, we further construct two core metrics from a global perspective. The **Mean Complete Uniformity (MCU)**, defined as

$$MCU = \frac{1}{m} \sum_{i=1}^m CU_i \quad (3)$$

measures the average balance level across all m attributes. The **Uniformity Coordination (UCO)**, given by

$$UCO = \frac{1}{m} \sum_{i=1}^m (CU_i - MCU)^2 \quad (4)$$

captures the dispersion of balance across attributes, reflecting their coordination consistency. Finally, by integrating **MCU** and **UCO**, we construct the **Penalized Global Balance Score**:

$$PGBS = MCU \times (1 - UCO) \quad (5)$$

PGBS simultaneously accounts for the uniformity of distributions within attributes and the coordination among attributes, resulting in a composite score bounded between $[0, 1]$. A higher score indicates a more desirable overall balance characteristic of the dataset. The balance of the dataset and the effectiveness of the Attribute Balancing Selection Algorithm are evaluated using the **PGBS** metric. The experimental results are shown in Tab. 25. **Our dataset achieves the highest balance score despite employing more extensive prompts. More significantly, this result was attained through a systematic and scalable algorithmic approach, eliminating the need for manual selection and its inherent biases.**

Table 22. Taxonomy of Predefined WordNet Synsets and Keywords (Major Content Aspect).

	WordNet Synsets	Key Phrases/Words
Actions	'travel.v.01', 'compete.v.01', 'act.v.01', 'manipulate.v.02', 'eat.v.01', 'chew.v.01', 'drink', 'move.v.03', 'move.v.02', 'change.v.01', 'make.v.03', 'make.v.01', 'run.v.01', 'crawl.v.01', 'stretch.v.01', 'bend.v.01', 'twist.v.01', 'balance.v.01', 'crouch.v.01', 'leap.v.01', 'analyze.v.01', 'memorize.v.01', 'predict.v.01', 'question.v.01', 'visualize.v.01', 'calculate.v.01', 'evaluate.v.01', 'meditate.v.01', 'negotiate.v.01', 'persuade.v.01', 'compliment.v.01', 'criticize.v.01', 'greet.v.01', 'introduce.v.01', 'console.v.01', 'debate.v.01', 'fold.v.01', 'sweep.v.01', 'organize.v.01', 'repair.v.01', 'pack.v.01', 'iron.v.01', 'polish.v.01', 'recycle.v.01', 'sketch.v.01', 'sculpt.v.01', 'compose.v.01', 'edit.v.01', 'design.v.01', 'photograph.v.01', 'choreograph.v.01', 'improvise.v.01', 'dribble.v.01', 'tackle.v.01', 'volley.v.01', 'parry.v.01', 'jog.v.01', 'sprint.v.01', 'strike.v.01', 'dive.v.01', 'diagnose.v.01', 'program.v.01', 'audit.v.01', 'negotiate.v.01', 'lecture.v.01', 'inspect.v.01', 'engineer.v.01', 'curate.v.01', 'scroll.v.01', 'upload.v.01', 'debug.v.01', 'encrypt.v.01', 'stream.v.01', 'click.v.01', 'zoom.v.01', 'render.v.01', 'sigh.v.01', 'glare.v.01', 'whisper.v.01', 'shrug.v.01', 'groan.v.01', 'cheer.v.01', 'murmur.v.01', 'giggle.v.01', 'harvest.v.01', 'dig.v.01', 'plant.v.01', 'prune.v.01', 'melt.v.01', 'evaporate.v.01', 'erode.v.01', 'ignite.v.01'	None
Kinetic Motions	'rotate.v.01', 'slide.v.01', 'pivot.v.01', 'vibrate.v.01', 'compress.v.01', 'expand.v.01', 'lock.v.01', 'unlock.v.01', 'accelerate.v.01', 'decelerate.v.01', 'lubricate.v.01', 'grind.v.01', 'weld.v.01', 'clamp.v.01', 'press.v.01', 'release.v.01', 'transmit.v.01', 'absorb.v.01', 'ignite.v.01', 'discharge.v.01', 'charge.v.01', 'heat.v.01', 'cool.v.01', 'filter.v.01', 'calibrate.v.01', 'adjust.v.01', 'trigger.v.01', 'activate.v.01', 'deactivate.v.01', 'sense.v.01', 'measure.v.01', 'align.v.01', 'cut.v.01', 'drill.v.01', 'screw.v.01', 'bolt.v.01', 'pump.v.01', 'vent.v.01', 'inject.v.01', 'extrude.v.01', 'brake.v.01', 'clutch.v.01', 'gear.v.01', 'lever.v.01', 'hydrate.v.01', 'insulate.v.01', 'conduct.v.01', 'resonate.v.01'	'rotate', 'move', 'bounce', 'spin', 'sway', 'flythrough', 'fly', 'panning', 'drone', 'run', 'walk', 'drive', 'zoom', 'chase', 'swim', 'movement', 'fall', 'rise', 'sliding video', 'sliding camera', 'sliding shot', 'forward', 'backward', 'leftward', 'rightward', 'upward', 'downward'
Fluid Motions	'body_of_water.n.01', 'fluid.n.01', 'fluid.n.02', 'atmospheric_phenomenon.n.01', 'deformation.n.02'	'fountain', 'float', 'firework', 'fire', 'cloud', 'clouds', 'candle', 'smoke', 'wave', 'inflate', 'melt', 'shrink', 'ripple'
Light Change	'burning.n.01', 'light.n.01', 'light.n.02', 'light.n.04', 'light.n.07', 'light.n.09'	'eclipse', 'sunset', 'sunrise', 'firework', 'fire', 'sunbeam', 'sun ray', 'sunshine', 'sunny', 'burn', 'shine', 'luminous', 'glow', 'explode', 'milky', 'galaxy', 'flash', 'sparkle', 'neon', 'reflection', 'bright', 'candle', 'light'

(a) Temporal categories under the "major content" aspect.

	WordNet Synsets	Key Phrases/Words
People	'person.n.01', 'people.n.01'	'he', 'she', 'men', 'team'
Animals	'animal.n.01'	None
Vehicles	'vehicle.n.01'	'drone'
Artifacts	'artifact.n.01'	None
Buildings & Infrastructures	'building.n.01', 'structure.n.01'	'building', 'cityscape', 'town', 'city'
Scenery & Natural Object	'natural_object.n.01', 'body_of_water.n.01', 'geological_formation.n.01', 'atmospheric_phenomenon.n.01', 'atmosphere.n.05'	'mountainous', 'fire', 'firework', 'solar eclipse', 'water current', 'water drop', 'cloud', 'desert'
Plants	'plant.n.02', 'vegetation.n.01'	None
Food & Beverage	'food.n.01', 'food.n.02'	'mushroom'
Illustrations	'shape.n.02', 'symbol.n.01'	'pattern', 'abstract', 'pattern', 'particle', 'gradient', 'loop', 'graphic'

(b) Spatial categories under the "major content" aspect.

Table 23. Taxonomy of Predefined WordNet Synsets and Keywords (Attribute Control Aspect).

	WordNet Synsets	Key Phrases/Words
Color	'color.n.01'	'white'
Camera View	"perspective.n.01", "angle.n.01", "view.n.01", "shot.n.01", "panorama.n.01", "closeup.n.01", "wideangle.n.01", "telephoto.n.01", "macro.n.01", "aerial.n.01", "birdview.n.01", "fisheye.n.01", "tilt.n.01", "pan.n.01", "zoom.n.01", "focus.n.01", "depth.n.01", "frame.n.01", "composition.n.01", "horizon.n.01", "vertical.n.01", "horizontal.n.01", "aspect.n.01"	'macro shot', 'medium shot', 'wide shot', 'close up', 'close-up', 'close view', 'close shot', 'front view', 'front-facing', 'front facing', 'backside view', 'backside shot', 'profile view', 'profile shot', 'side view', 'side shot', 'top view', 'top-down view', 'top down view', 'overhead view', 'overhead shot', 'bottom view', 'bottom shot', 'low angle', 'high angle', 'aerial', 'drone view', 'bird's eye view', 'first person', 'first-person', '1st person', 'third person', 'third-person', '3rd person'
Quantity	'integer.n.01'	'a'

(a) Spatial categories under the "attribute control" aspect.

	WordNet Synsets	Key Phrases/Words
Speed	"velocity.n.01", "speed.n.01", "acceleration.n.01", "deceleration.n.01", "momentum.n.01", "rapidity.n.01", "swiftness.n.01", "haste.n.01", "pace.n.01", "tempo.n.01", "vector.n.01", "terminal.n.01", "escape.n.01", "angular.n.01", "linear.n.01", "relative.n.01", "initial.n.01", "final.n.01", "constant.n.01", "variable.n.01", "light.n.01", "sound.n.01"	'slow', 'fast', 'slowly', 'fastly', 'timelapse', 'timelapse', 'time-lapse', 'stop motion', 'velocity', 'speed', 'acceleration', 'deceleration', 'momentum', 'rapidity', 'swiftness', 'haste', 'pace', 'tempo', 'velocity vector', 'terminal velocity', 'escape velocity', 'angular velocity', 'linear velocity', 'relative velocity', 'initial velocity', 'final velocity', 'constant speed', 'variable speed', 'speed of light', 'speed of sound'
Motion Direction	"north.n.01", "south.n.01", "east.n.01", "west.n.01", "forward.n.01", "backward.n.01", "left.n.01", "right.n.01", "up.n.01", "down.n.01", "vertical.n.01", "horizontal.n.01", "radial.n.01", "tangential.n.01", "axial.n.01", "diagonal.n.01", "lateral.n.01", "longitudinal.n.01", "reverse.n.01", "clockwise.n.01", "counterclockwise.n.01", "oblique.n.01", "azimuth.n.01", "bearing.n.01", "orientation.n.01", "trajectory.n.01", "normal.n.01", "parallel.n.01", "perpendicular.n.01"	'ahead', 'anticlockwise', 'away from', 'clockwise', 'counterclockwise', 'downward', 'eastbound', 'northbound', 'southbound', 'westbound', 'homeward', 'leftwards', 'rightwards', 'upward', 'left', 'right', 'forward', 'backward', 'toward', 'out of', 'approach', 'leave', 'against to', 'lift', 'opposite direction'
Event Order	"sequence.n.01", "order.n.01", "phase.n.01", "step.n.01", "stage.n.01", "priority.n.01", "predecessor.n.01", "successor.n.01", "time.n.01", "past.n.01", "future.n.01", "present.n.01", "duration.n.01", "interval.n.01", "epoch.n.01", "era.n.01", "period.n.01", "schedule.n.01", "calendar.n.01", "chronology.n.01", "timing.n.01", "frequency.n.01", "cycle.n.01", "simultaneous.n.01", "anachronism.n.01"	'and then'

(b) Temporal categories under the "attribute control" aspect.

Table 26. Impact of Video Compression Formats on Detector Performance.

Model	Training Scheme	Accuracy (%)		Gap
		H.264 Test	MPEG-4 Test	
I3D	Mixed	100.00	100.00	0.00
	Unified	100.00	100.00	0.00
X3D	Mixed	99.93	99.93	0.00
	Unified	99.97	100.00	0.033
CNN_Spot	Mixed	99.97	99.97	0.00
	Unified	99.97	99.97	0.00
UnivFD	Mixed	99.83	99.83	0.00
	Unified	99.83	99.83	0.00

- Mixed: Real videos (MPEG-4), Generated videos (H.264);
- Unified: All videos compressed with H.264.

8.6. Discussion on Video Encoding Methods

In early explorations of generated image detection [56, 80, 103], an interesting phenomenon was observed: real videos were typically stored in PNG format, while generated images were saved as JPEG. Models [80] trained on such datasets could easily distinguish between positive and negative samples, achieving high accuracy. This unintentional format discrepancy led the models to learn the compression differences between PNG and JPEG, rather than the intrinsic

distinctions between authentic and generated content.

Motivated by this finding, we investigated whether a similar effect would occur when real videos are compressed using MPEG-4 Part 2 and generated samples are compressed with H.264. As shown in Tab. 26, we selected several common network architectures and trained them under two conditions: one where real videos were compressed with MPEG-4 Part 2 and generated samples with H.264, and another where all videos were uniformly compressed with H.264. The models were then evaluated on real videos compressed with either MPEG-4 Part 2 or H.264. From the perspective of accuracy metrics, the impact of this compression discrepancy is negligible. Only the X3D model, trained on the mixed-compression dataset (real: MPEG-4 Part 2, generated: H.264), exhibited any noticeable effect. Furthermore, although we observed some variations in output logits Fig. 24, these changes had almost no influence on the final accuracy. Nevertheless, to mitigate potential confounding factors, all videos in AIGVDBench were uniformly compressed using H.264.

8.7. Training details

Training Protocol for Non-VLM Models: To maximize the inherent potential of each individual model (excluding VLMs), we deliberately refrained from standardizing key hyperparameters such as learning rates and data augmenta-

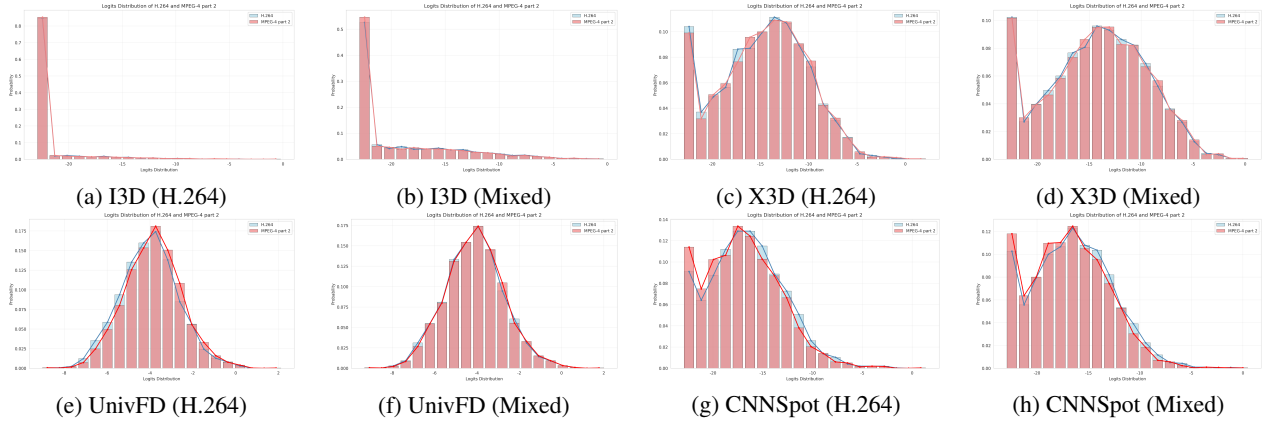
Table 24. Examples of Prompt Classification Results. Misclassified properties are denoted by red text.

Prompt	Spatial Content	Temporal Content	Spatial Attribute	Temporal Attribute
two men are engaged in a conversation in a car showroom. The older man, dressed in a black polo shirt, is speaking to the younger man, who is wearing a black t-shirt. The showroom is filled with various cars, including a Bugatti, and posters adorn the walls. The men are standing in front of a car, suggesting that they might be discussing the vehicle. . .	People Artifacts Buildings&Infrastructures Vehicles	Actions	Color Quantity	None
a close-up of a doll with curly brown hair and large, expressive eyes. The doll is wearing a denim dress with a floral design on the front. In the background, there are cans of pumpkin pie filling, suggesting a setting related to baking or cooking. . .	Artifacts Scenery&Natural Object Plants Food&Beverage	Actions Light Change	Color Camera View Quantity	None
three young men standing on a race track, smiling and posing for the camera. They are dressed in casual clothing, with one wearing a tie-dye shirt, another in a tank top, and the third in a black t-shirt with a cat design. . .	People Animals Artifacts Buildings&Infrastructures Scenery&Natural Object	Actions Fluid Motions	Color Quantity	None
a white Mustang GT500 parked in a parking lot. The car is adorned with blue stripes on the hood and side, and it features a black grille and black wheels. The car is parked in a spot with a white line marking the boundary . . .	People Animals Vehicles Artifacts Illustrations	Actions Kinetic Motions Light Change	Color Camera View Quantity	None
a young man standing next to a red Volkswagen Beetle car on a road surrounded by trees. The man is dressed in a casual blue shirt and jeans, and he has his hands clasped together. The car is parked on the side of the road, and it has a shiny red finish. . .	People Animals Artifacts Scenery&NaturalObject Plants	Actions Fluid Motions LightChange	Color Camera View Quantity	None
a white SUV driving down a snowy road. The vehicle is equipped with a roof rack and has a distinctive orange stripe running along its side. The road is lined with trees, and the snow-covered landscape suggests a cold, winter day . . .	People Vehicles Artifacts Buildings&Infrastructures Scenery&NaturalObject Plants	Actions Kinetic Motions	Color Camera View Quantity	Speed Motion Direction

tion schemes. This decision was motivated by the fact that such parameters are intrinsically tied to each model’s ar-

chitecture, and arbitrary unification could compromise their original performance. Furthermore, given the large num-

Figure 24. Logit distribution comparison across different models under two training schemes.



Note: H.264 = models trained with unified H.264 compression; Mixed = models trained with mixed compression (real videos: MPEG-4 Part 2, generated samples: H.264). The figure shows feature distributions when tested on datasets containing both compression formats. In all subfigures, the light coral curves represent test results on MPEG-4 Part 2 compressed videos, while the light blue curves represent test results on H.264 compressed videos.

Table 25. Balance Comparison with Prior Datasets($\alpha = 2.0$).

Metric	AIGVDBench	GVF	GenVidBench
Complete Uniformity (CU) by Attribute			
Spatial_Content	0.9592	0.7818	0.8495
Temporal_Content	0.8785	0.6866	0.7677
Spatial_Attribute	0.9942	0.4069	0.7492
Temporal_Attribute	1.0000	0.8762	0.9966
Global Metrics			
MCU	0.9580	0.6879	0.8408
UCO	0.0023	0.0308	0.0095
PGBS	0.9557	0.6667	0.8328
Prompt Length Analysis			
Avg. Prompt Length	552.39	46.29	88.69

ber of models involved, enforcing a common parameter set was deemed impractical. To mitigate overfitting, we set an upper limit of 50 training epochs for all models, ensuring that each had adequately converged before termination. For models originally proposed with fewer than 50 epochs in their respective papers, we adhered to their prescribed training length. In terms of input resolution, most models were trained and evaluated using 256x256 images, while Vision Transformer (ViT)-based models utilized 224x224 inputs. For video classification models, we adopted the MMac-tion2 [59] and followed the recommended hyperparameter settings. For generated image detection, early-stage methods were implemented based on the AIGCDetectBenchmark [99] repository, while other models were trained using officially released code. Across all experiments, the only modification made was to the data input pipeline; all other training configurations were preserved as provided in the original implementations.

Table 27. Prompt template assignment for VLM models.

Prompt Template Group	Models
Common Template	Qwen2.5-VL-3B-Instruct
	Qwen2.5-VL-7B-Instruct
	Qwen2.5-VL-32B-Instruct
	Kimi-VL-A3B-Instruct
	FastVLM-Apple-7B
	InternVL-8B
Individual Templates	LLaVA-1.5-7B
	Emu3-Stage1
	DeepSeek-VL-7B
	DeepSeek-VL2-Small (DeepSeek-VL2)

VLM Models: It must be acknowledged that Vision-Language Models (VLMs) have demonstrated remarkable potential across a range of tasks. In our experimental setup, we adopted a prompt-based zero-shot approach to objectively evaluate VLMs’ capabilities in generated video detection. Due to variations in prompt-following styles across different models, we adjusted the prompts for models that initially exhibited poor classification performance to ensure a fair comparison of their “basic classification capabilities”. The specific prompt templates are illustrated in Figs. 25 to 28. Fig. 25 presents the common prompt template employed for multiple VLM models, while Figs. 26 to 28 display the individually tailored prompt templates designed for Emu3-Stage1, DeepSeek-VL-7B, and DeepSeek-VL2-Small respectively. For Vision-Language Models (VLMs) lacking temporal reasoning capabilities, we implement a

```

Common Prompt Template

## Critical Instruction
You MUST respond with EXACTLY ONE WORD only.

## Task
Analyze this video frame for deepfake detection. Determine if the content is real or AI-generated.

## Detection Criteria
Look for the following features:



- Visual artifacts
- Inconsistencies
- Unnatural movements
- Authentication markers


## Response Format



- If content is authentic/real: Answer 'REAL'
- If content is AI-generated/deepfake: Answer 'FAKE'


## Strict Rules



- Output ONLY the classification word
- No explanations, descriptions, or other text allowed
- Any response other than 'REAL' or 'FAKE' will be considered invalid

```

Figure 25. Common prompt template for AI-generated video detection, used with multiple VLM models.

```

Emu3-Stage1 Prompt Template

## Role
You are analyzing video frames for deepfake detection.

## Task
Look carefully at these frames and determine if they show real content or AI-generated/deepfake content.

## Detection Criteria
Consider the following features:



- Visual artifacts
- Inconsistencies
- Unnatural movements
- Authenticity markers


## Examples



- If content is authentic: Answer 'REAL'
- If content is AI-generated/deepfake: Answer 'FAKE'


## Important Rules



- You must respond with ONLY one word - either 'REAL' or 'FAKE'
- Do not provide any explanation, description, or other text
- Just the single classification word

```

Figure 26. Individual prompt template for AI-generated video detection with Emu3-Stage1 model.

```

DeepSeek-VL-7B Prompt Template

## <Image_placeholder>

## Role
You are a forensic deepfake detector. Inspect the provided video frame carefully.

## Task
Decide whether the frame is REAL (authentic) or FAKE (AI-generated / manipulated).

## Focus Areas
Analyze the following aspects:



- Artifacts
- Lighting inconsistencies
- Texture glitches
- Boundary errors
- Unnatural details


## Critical Instruction
Respond with EXACTLY ONE WORD: 'REAL' or 'FAKE'.

[CONTEXT PLACEHOLDER]
Classification:

```

Figure 27. Individual prompt template for AI-generated video detection with DeepSeek-VL-7B model.

```

DeepSeek-VL2 Prompt Template

## <Image>

## Role
You are a forensic deepfake detector. Some samples are suspicious.

## Task
Analyze this video frame for deepfake detection. Determine if the content is real or AI-generated.

## Detection Criteria
Look for the following features:



- Visual artifacts
- Inconsistencies
- Unnatural movements
- Authentication markers


## Critical Rule
Respond ONLY with one word: REAL or FAKE.

[CONTEXT PLACEHOLDER]
Classification:

```

Figure 28. Individual prompt template for AI-generated video detection with DeepSeek-VL2-Small model.

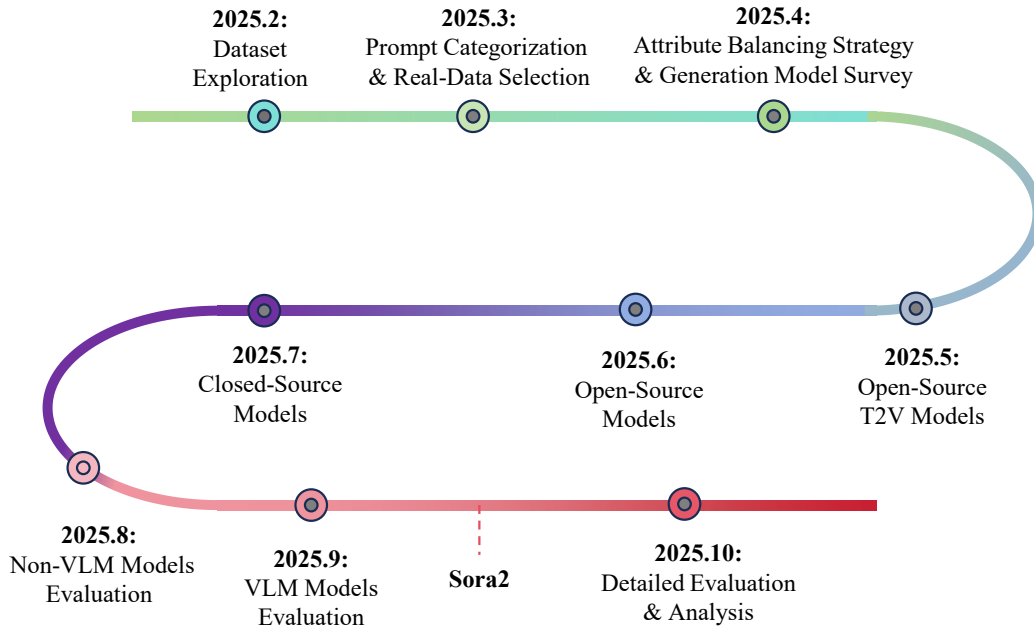


Figure 29. Timeline for AIGVDBench: Key Milestones and Development Phases from February to October 2025.

frame-based evaluation protocol: a video is labeled as "REAL" if at least one frame is classified as "REAL", and as "FAKE" only if all frames are consistently classified as "FAKE". In cases where the model fails to produce a valid output, the video is categorized as "No Answer".

Hardware Configuration: During the dataset construction phase, we utilized eight RTX 3090 GPUs and one 80GB A800 GPU. For the evaluation stage, all non-VLM models were assessed using four RTX 3090 GPUs, while the evaluation of VLM models was conducted on a single 80GB A800 GPU.

8.8. Timeline of AIGVDBench Construction

As illustrated in Fig. 29, the project followed a structured development process with key milestones achieved from February to October 2025.

9. Limitation

In this paper, we introduce AIGVDBench and conduct a series of explorations based on it. While substantial efforts and resources have been devoted to this work, several limitations should be acknowledged. From a data perspective, although we aimed to encompass a wide range of existing generative models, we were constrained by computational

resources and generation speed, which necessitated the use of smaller-parameter versions of some models. Additionally, due to the high cost associated with generating samples using closed-source models, we resorted to collecting existing videos rather than regenerating the entire test set from scratch. Furthermore, certain closed-source models, such as Sora 2 and Veo 3, were excluded due to regional access restrictions.

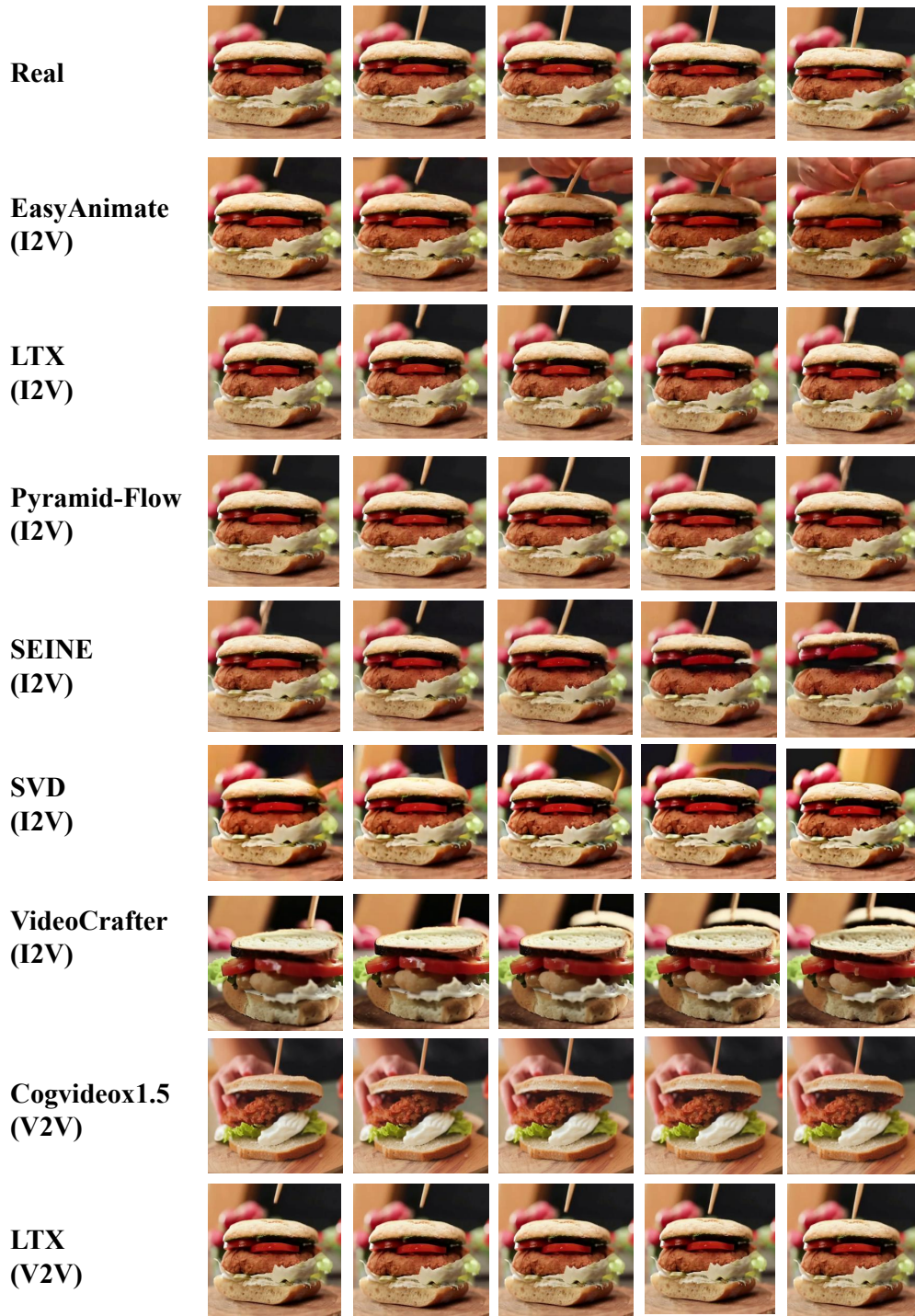
In terms of evaluation, while we incorporated several state-of-the-art models, storage and time limitations led to the exclusion of certain approaches, such as reconstruction-based methods. All analytical findings are derived from rigorously designed experiments; however, we cannot guarantee absolute correctness, as the dataset construction process may have introduced extraneous factors, despite our careful efforts to minimize such interference.

We hope that this work will facilitate and accelerate research in the field of generated video detection.

10. More examples of AIGVDBench

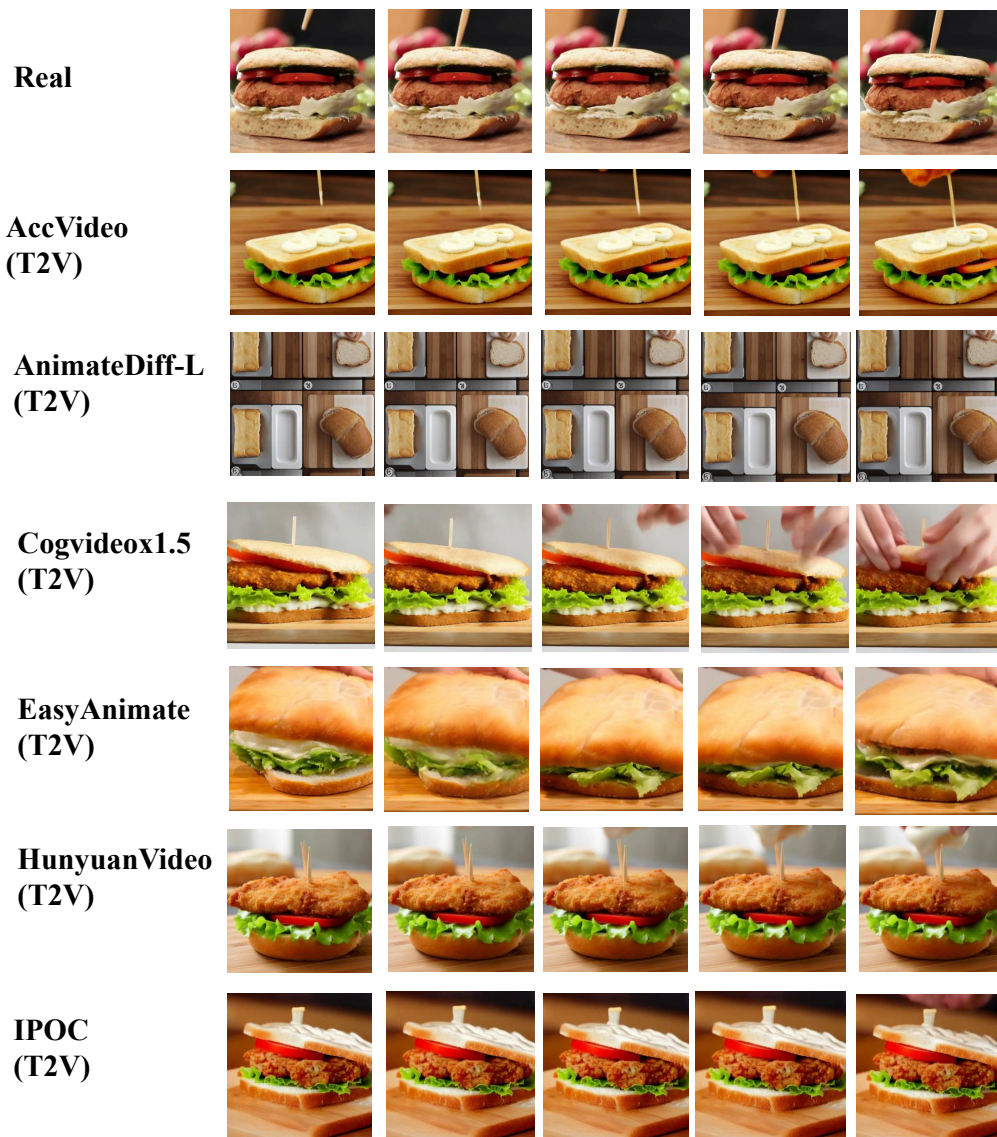
10.1. Video Samples from the AIGVDBench

Additional video examples are presented in Figs. 30 to 39.



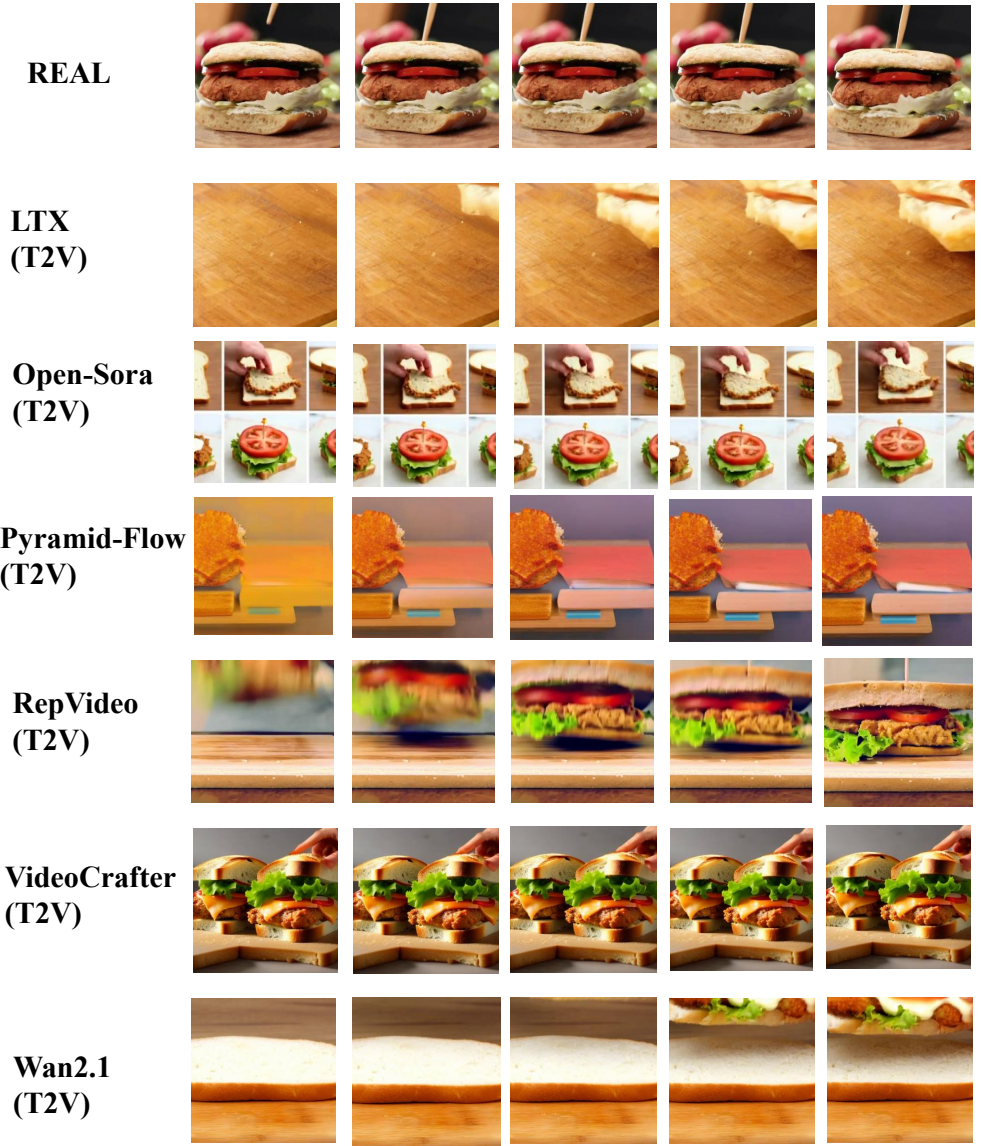
Prompts:
 The process of making a delicious sandwich. In the first frame, a wooden cutting board is placed on a table ... The style of the video is simple and straightforward, focusing on the sandwich-making process without any additional elements or distractions.

Figure 30. Video Visualization Results from Open-Source Models on AIGVDBench.



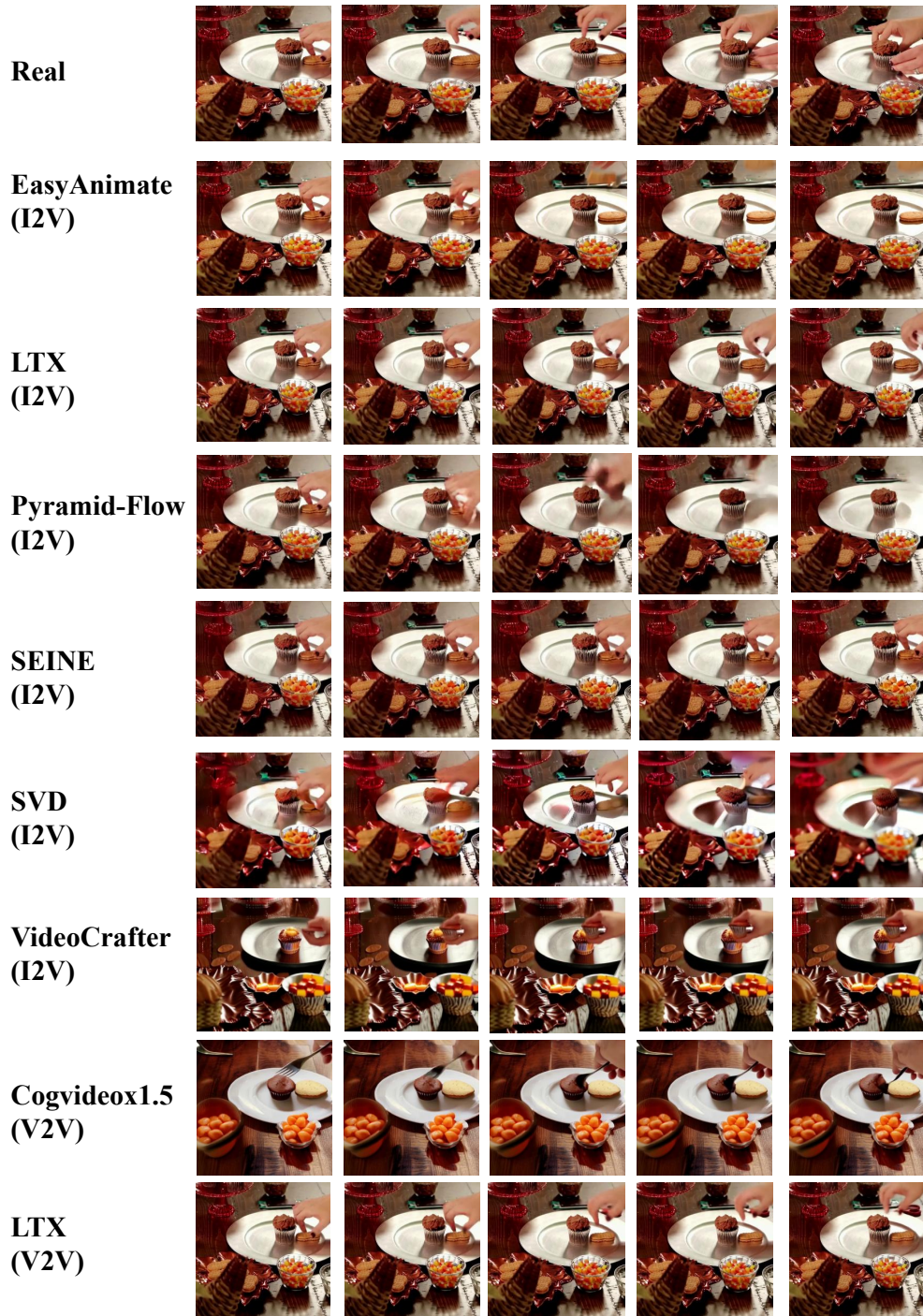
Prompts:
The process of making a delicious sandwich. In the first frame, a wooden cutting board is placed on a table ... The style of the video is simple and straightforward, focusing on the sandwich-making process without any additional elements or distractions.

Figure 31. Video Visualization Results from Open-Source Models on AIGVDBench.



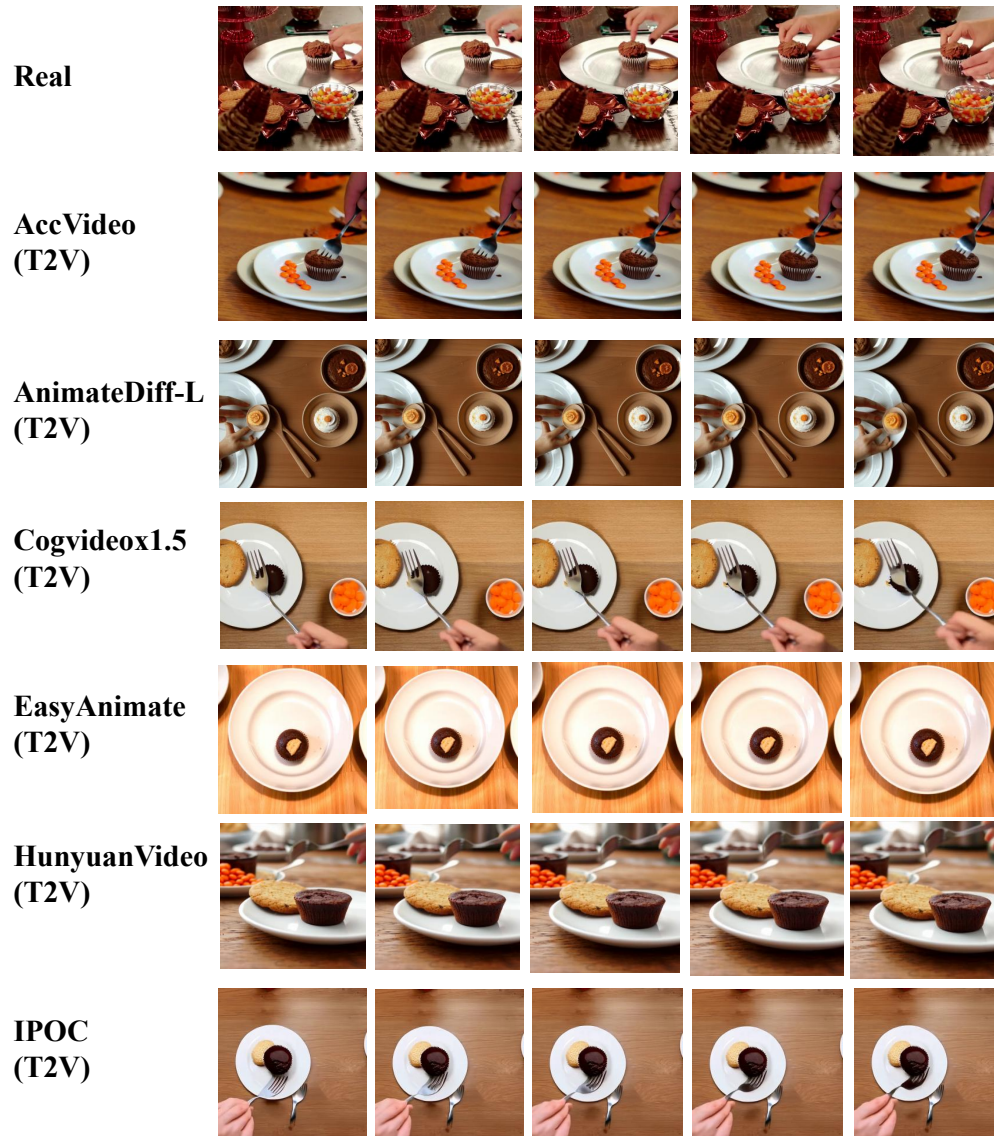
Prompts:
The process of making a delicious sandwich. In the first frame, a wooden cutting board is placed on a table ... The style of the video is simple and straightforward, focusing on the sandwich-making process without any additional elements or distractions.

Figure 32. Video Visualization Results from Open-Source Models on AIGVDBench.



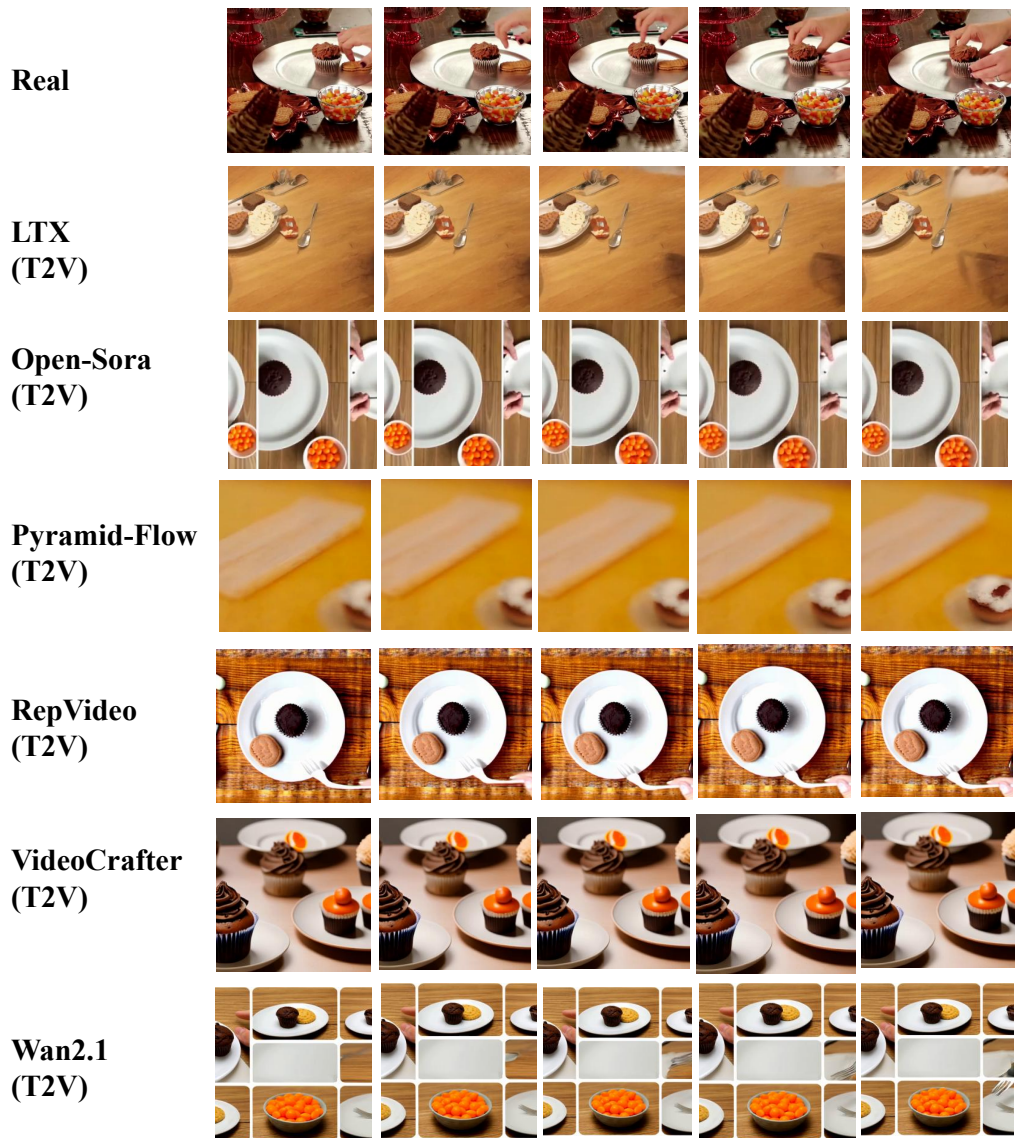
Prompts:
 A close-up of a table with various food items and utensils. The table is wooden with a visible grain pattern. ... everyday scene with a focus on the food and the act of eating. The lighting is bright and even, highlighting the textures and colors of the food items."

Figure 33. Video Visualization Results from Open-Source Models on AIGVDBench.



Prompts:
 A close-up of a table with various food items and utensils. The table is wooden with a visible grain pattern. ... everyday scene with a focus on the food and the act of eating. The lighting is bright and even, highlighting the textures and colors of the food items."

Figure 34. Video Visualization Results from Open-Source Models on AIGVDBench.



Prompts:
 A close-up of a table with various food items and utensils. The table is wooden with a visible grain pattern. ... everyday scene with a focus on the food and the act of eating. The lighting is bright and even, highlighting the textures and colors of the food items."

Figure 35. Video Visualization Results from Open-Source Models on AIGVDBench.

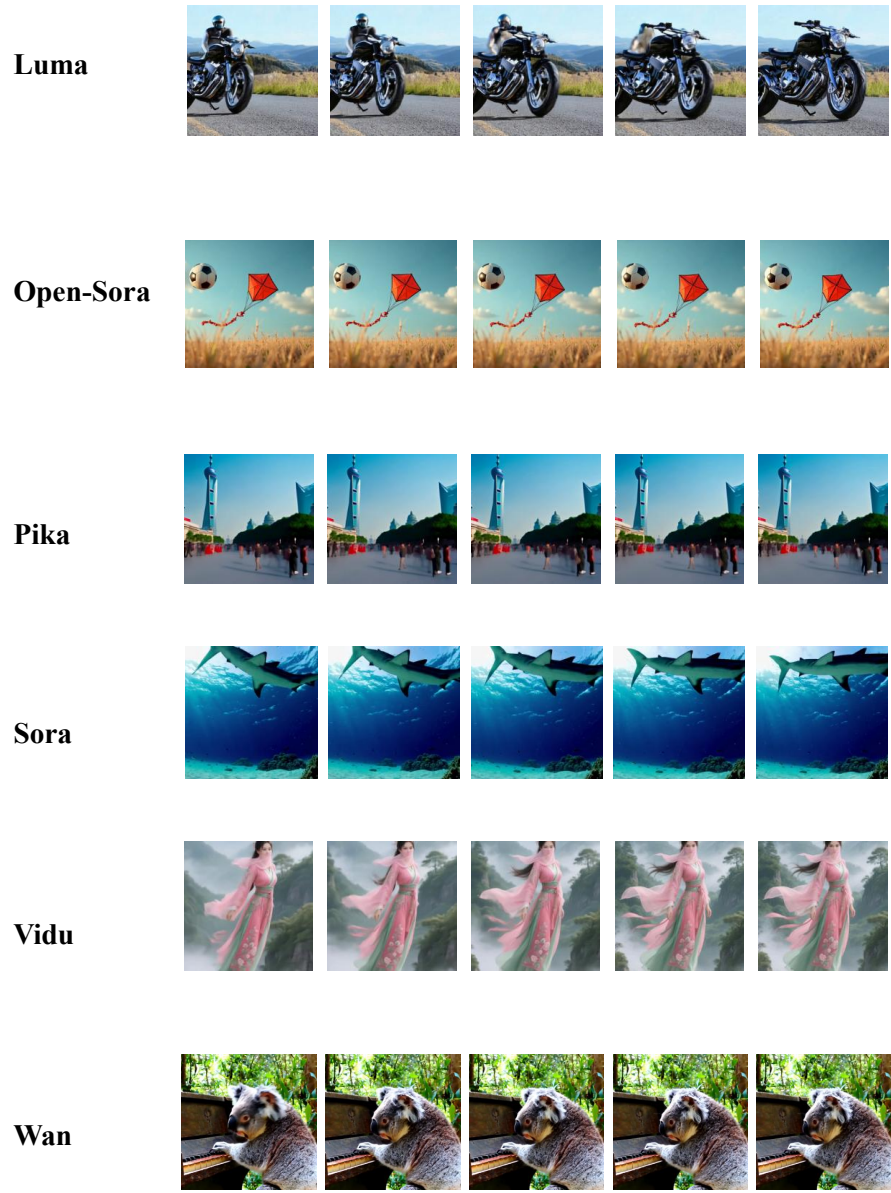
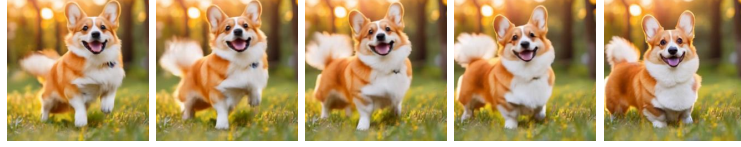
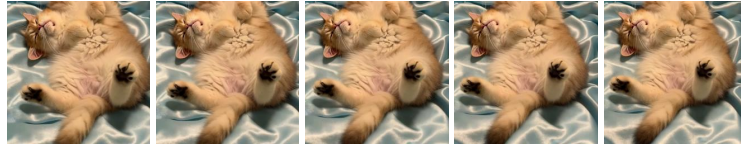


Figure 36. Video Visualization Results from Closed-Source Models on AIGVDBench.

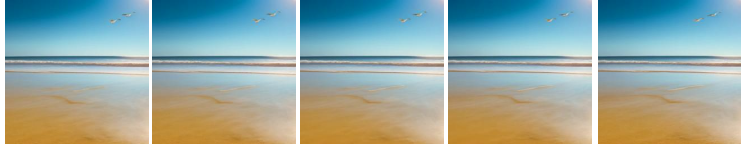
Luma



Open-Sora



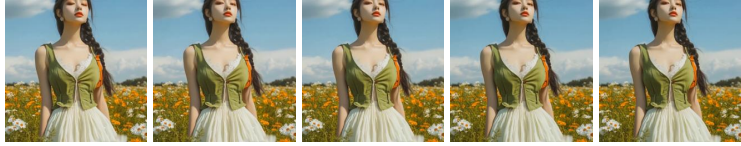
Pika



Sora



Vidu



Wan

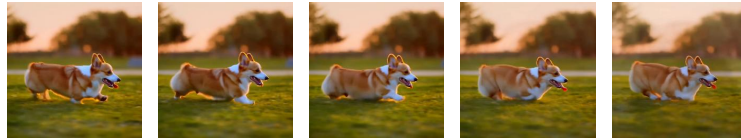


Figure 37. Video Visualization Results from Closed-Source Models on AIGVDBench.

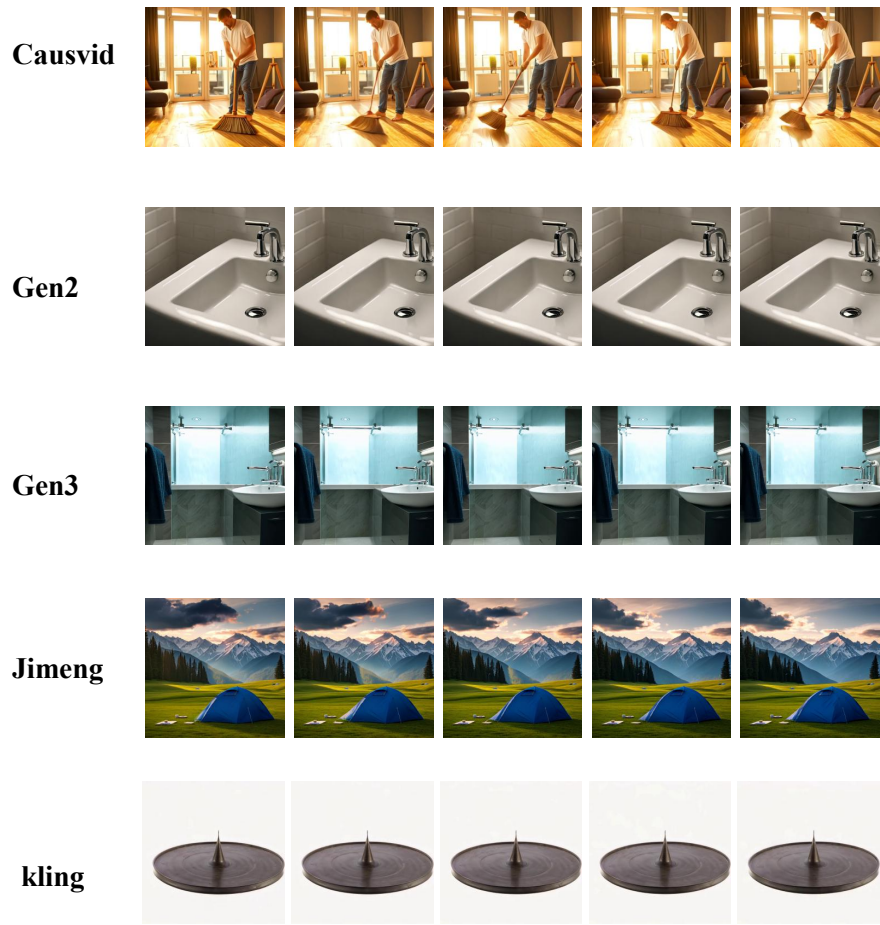


Figure 38. Video Visualization Results from Closed-Source Models on AIGVDBench.

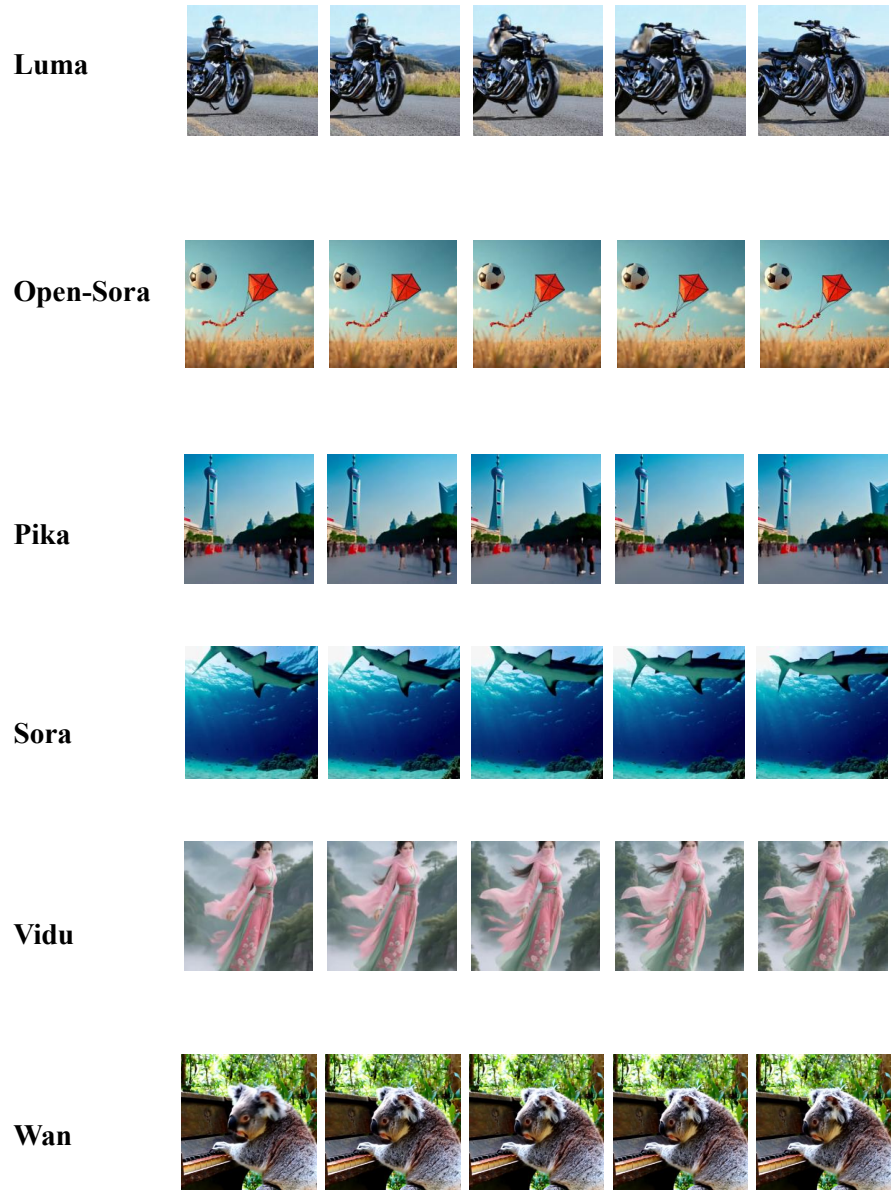
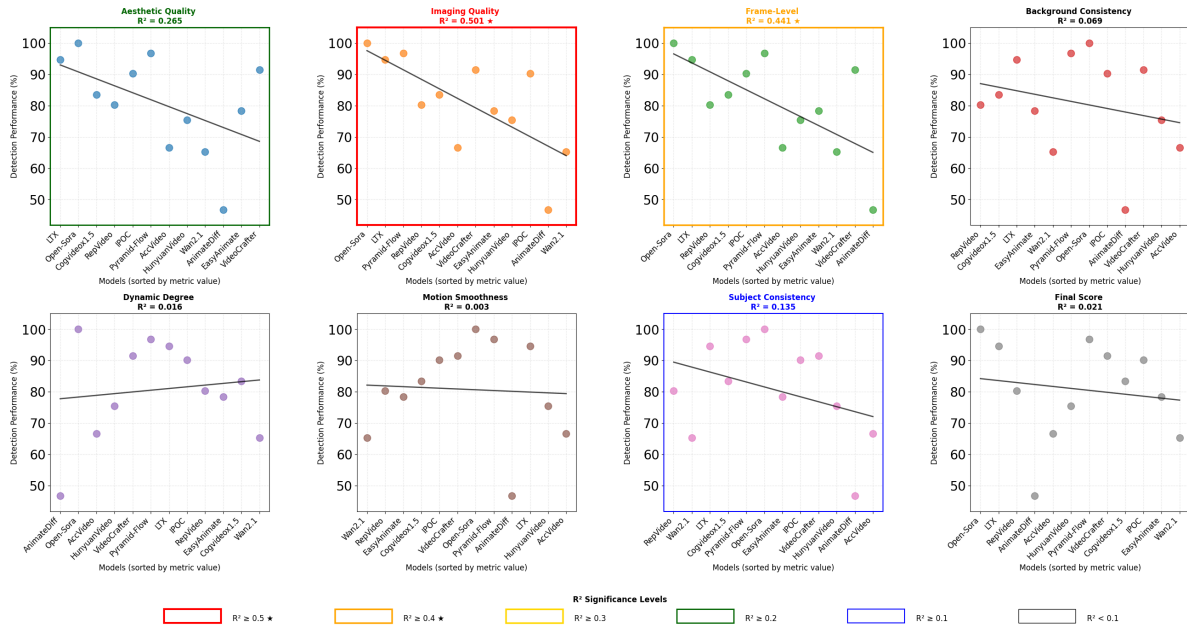


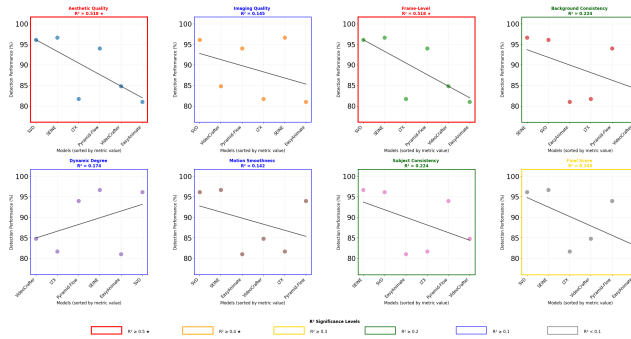
Figure 39. Video Visualization Results from Closed-Source Models on AIGVDBench.

10.2. More Experimental Results

The complete experimental results of the Analysis-2.1 Experiments A (Fig. 40- Fig. 59) and B (Fig. 62- Fig. 79) are presented in this section.



(a) Text-to-Video (T2V) task



(b) Image-to-Video (I2V) task

Figure 40. Correlations between the performance of I3D and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

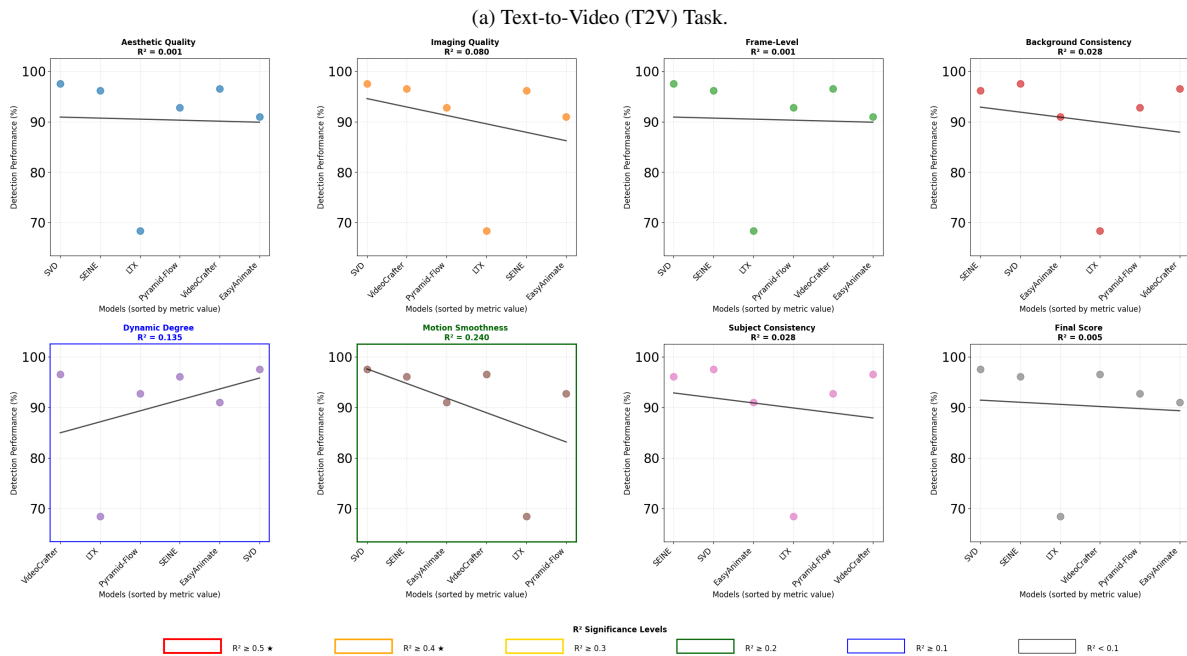
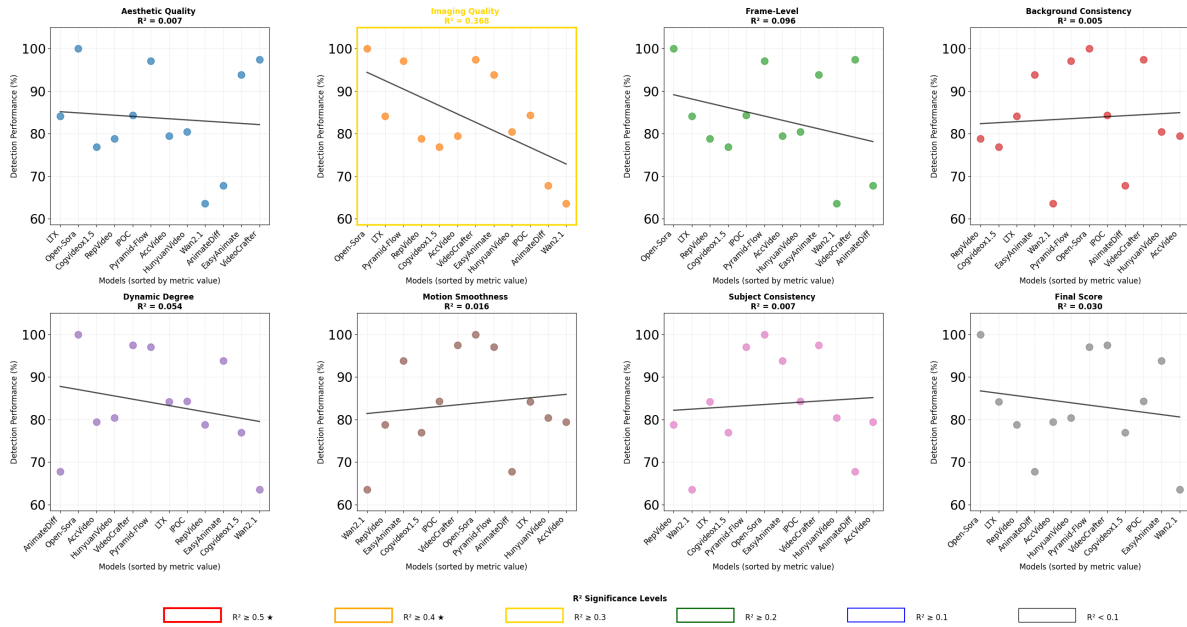
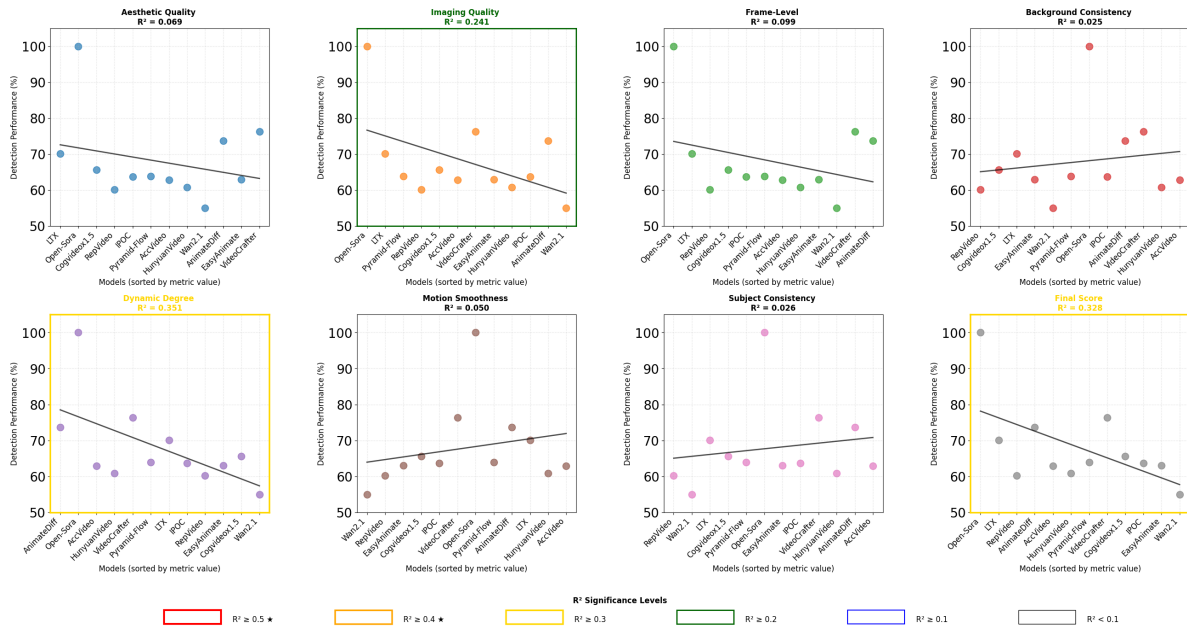
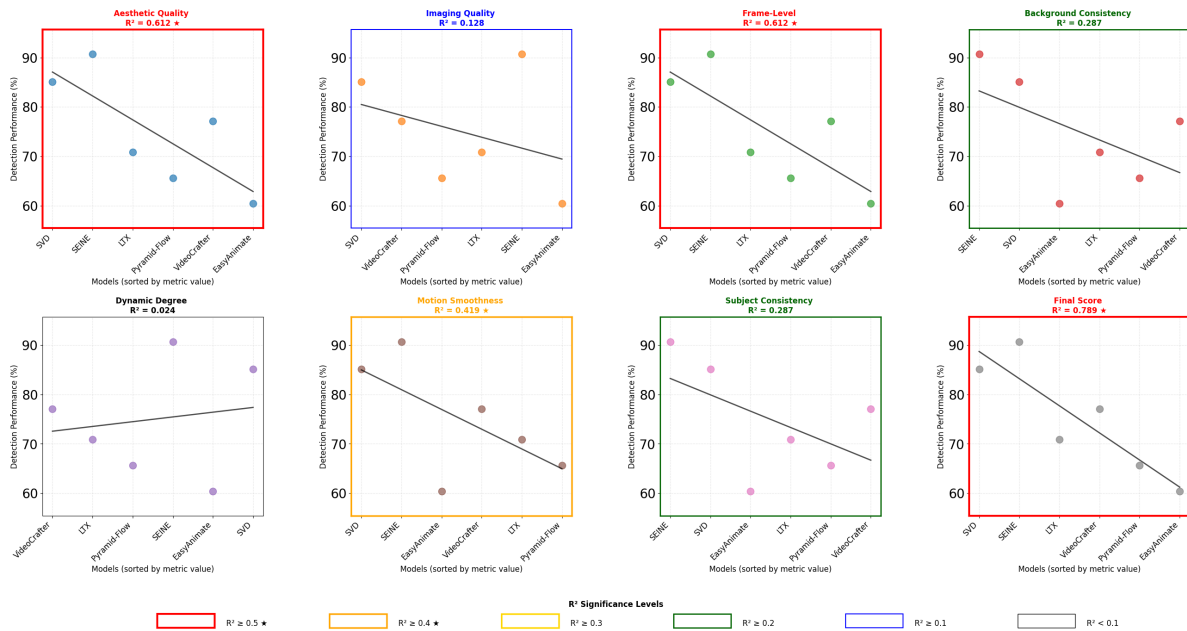


Figure 41. Correlations between the performance of X3D and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).



(a) Text-to-Video (T2V) Task.



(b) Image-to-Video (I2V) Task.

Figure 42. Correlations between the performance of VideoSwin and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

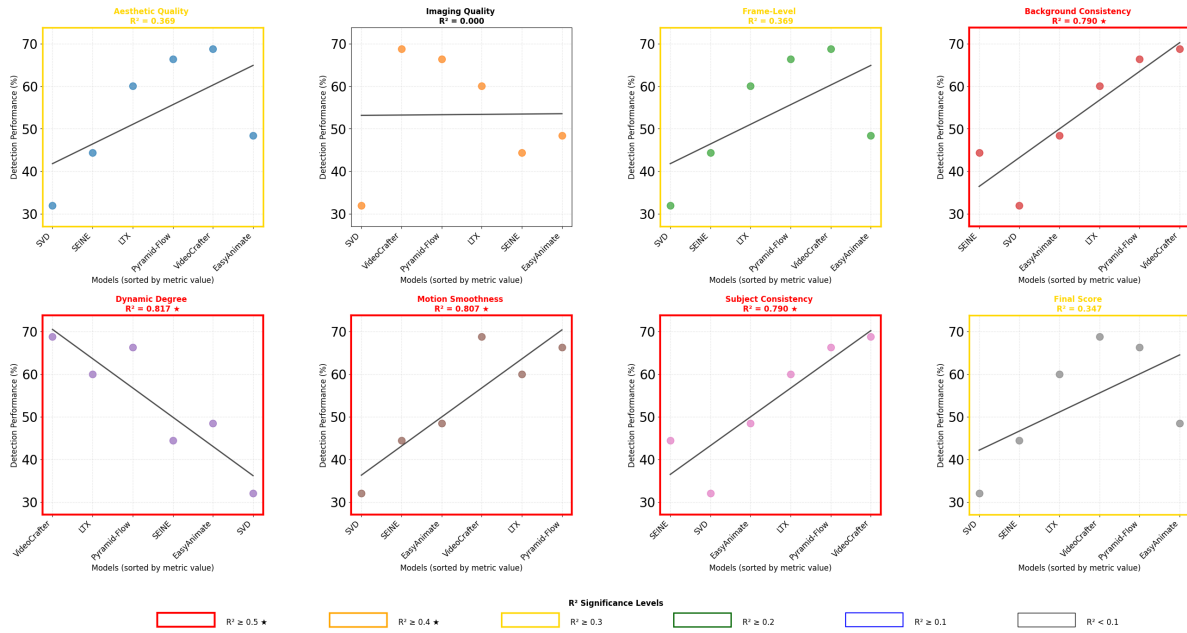
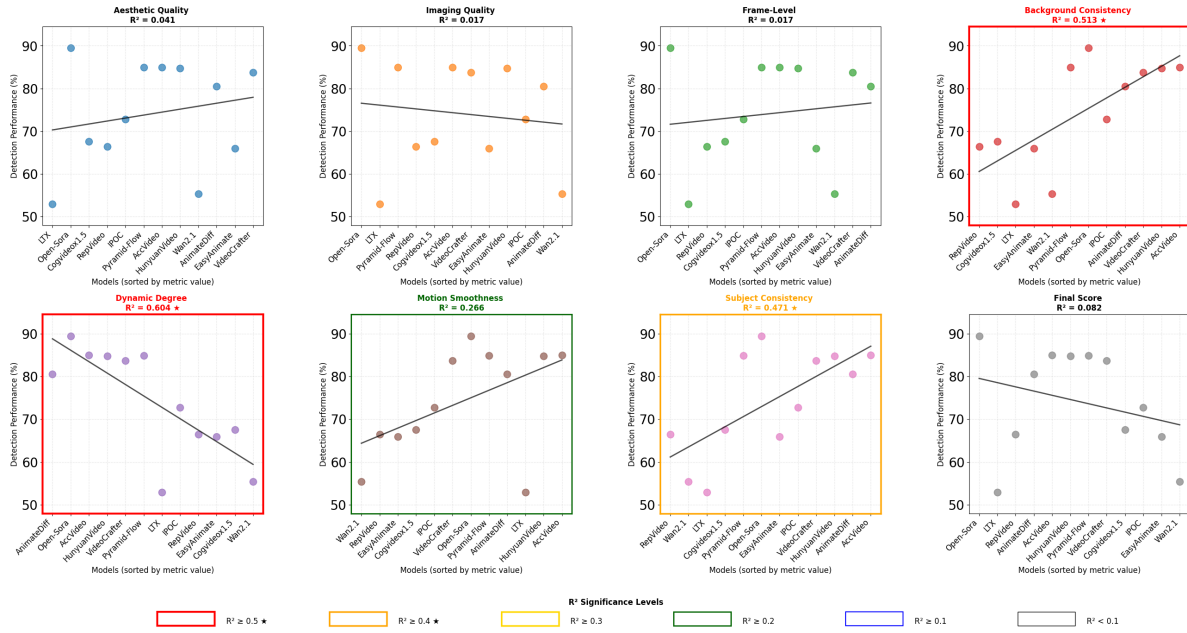
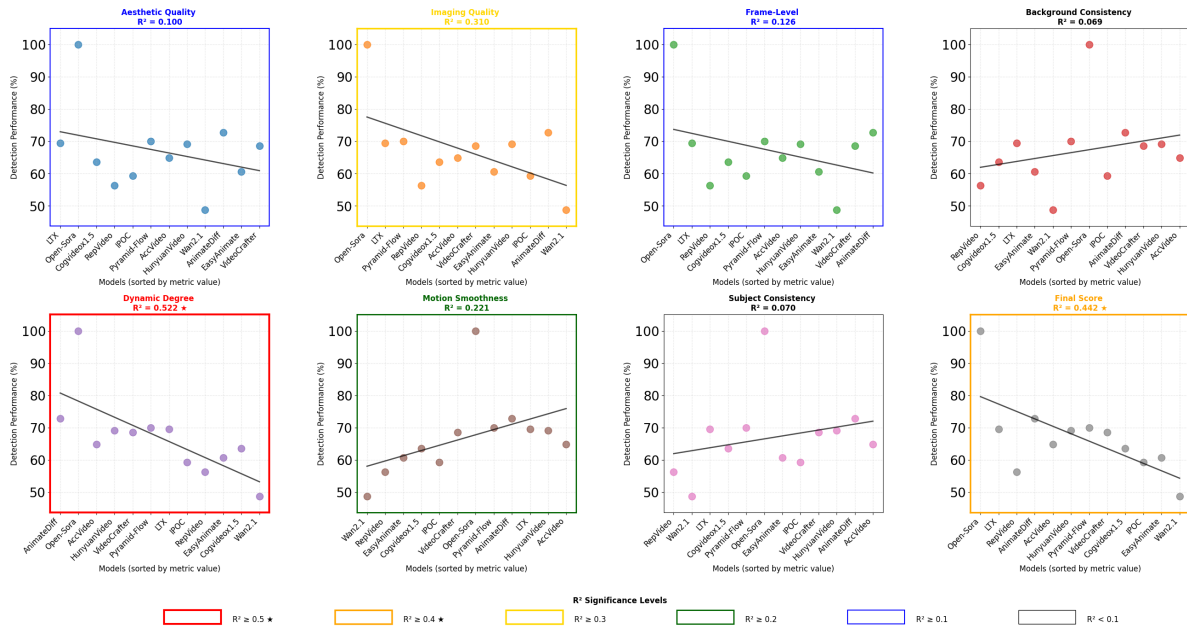
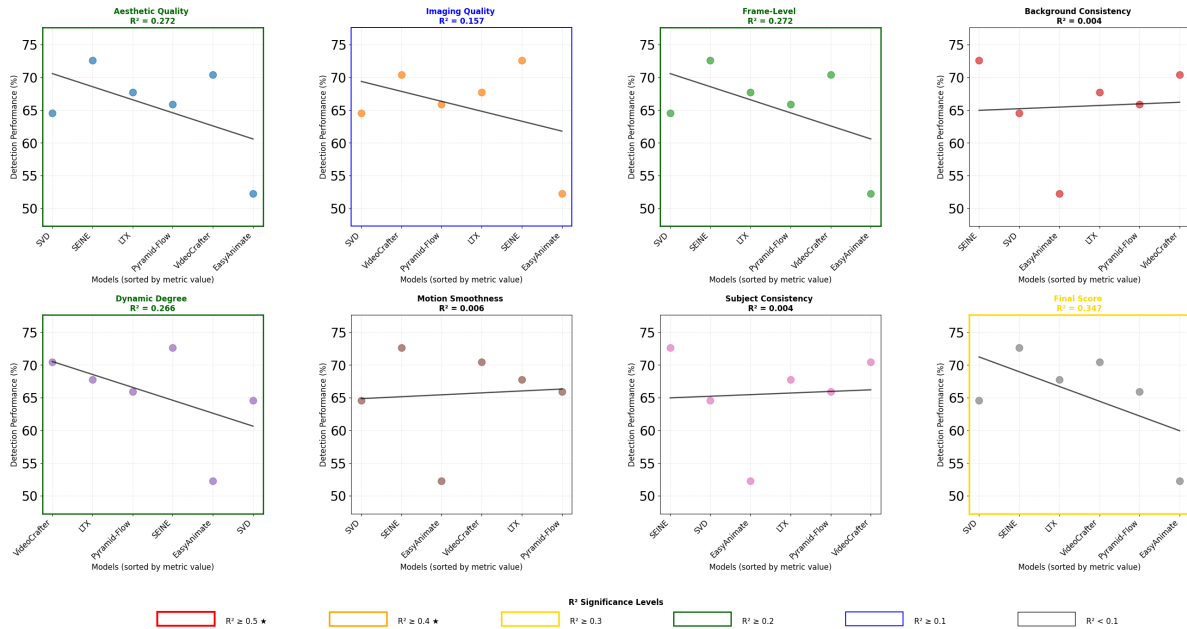


Figure 43. Correlations between the performance of UniFormer and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

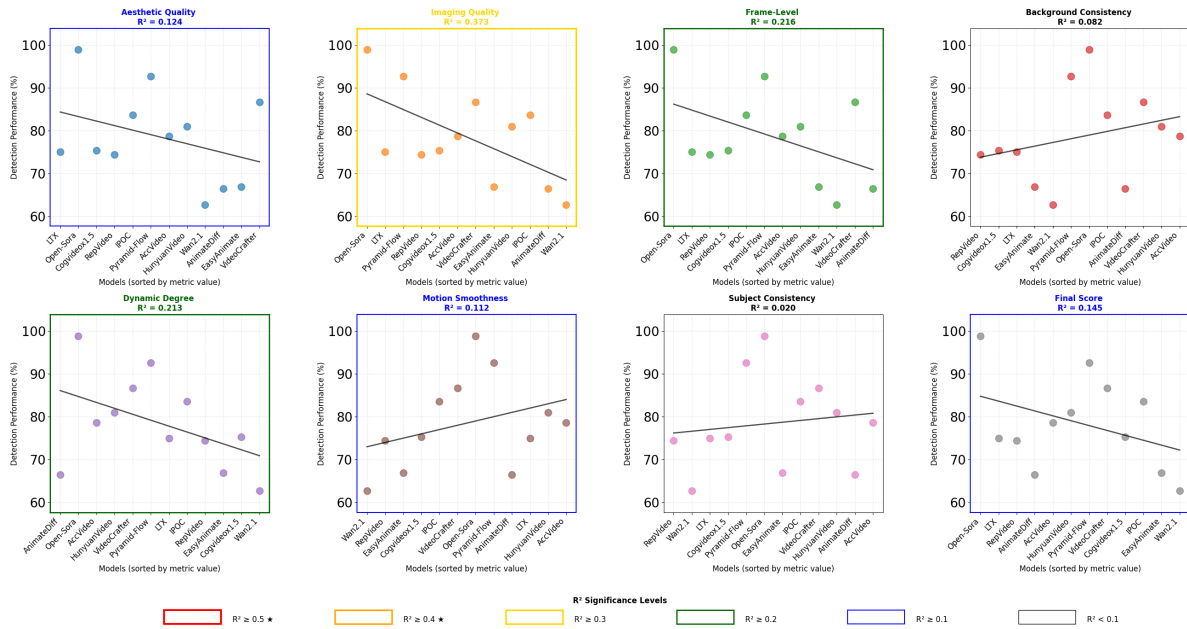


(a) Text-to-Video (T2V) Task.

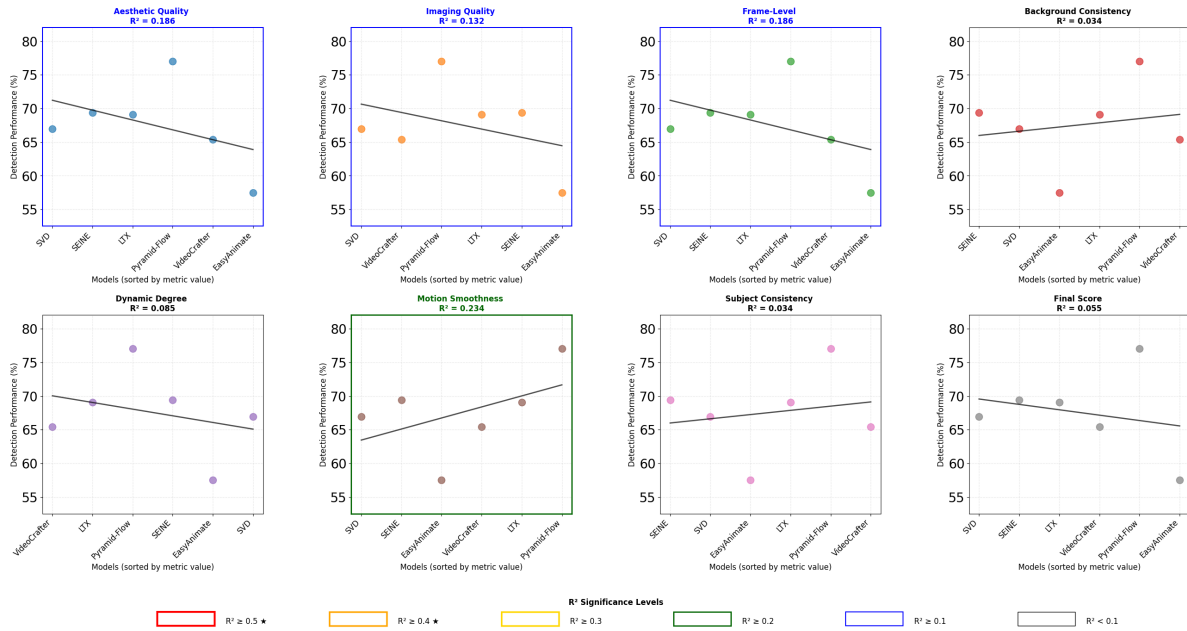


(b) Image-to-Video (I2V) Task.

Figure 44. Correlations between the performance of MVItv2 and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).



(a) Text-to-Video (T2V) Task.



(b) Image-to-Video (I2V) Task.

Figure 45. Correlations between the performance of VideoMAE and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

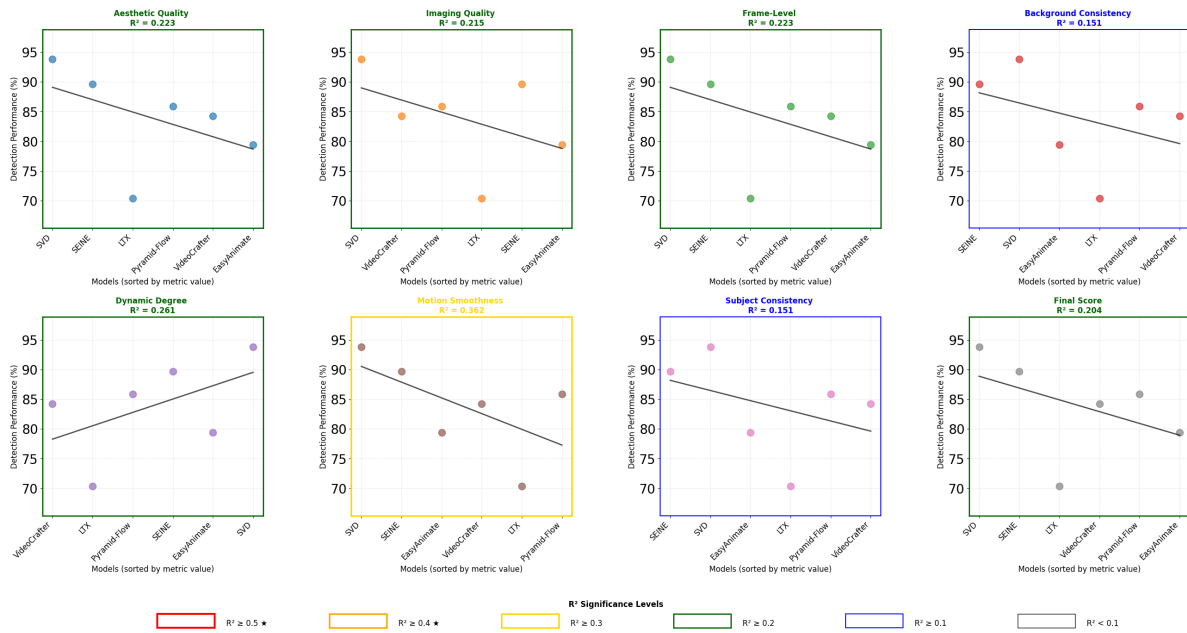
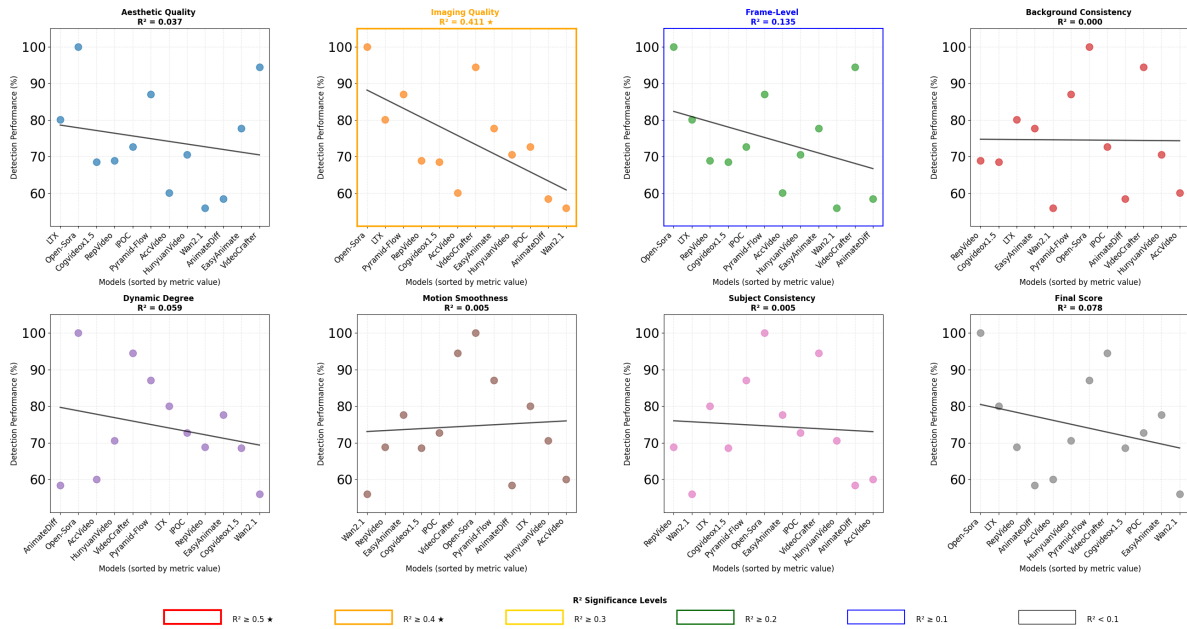
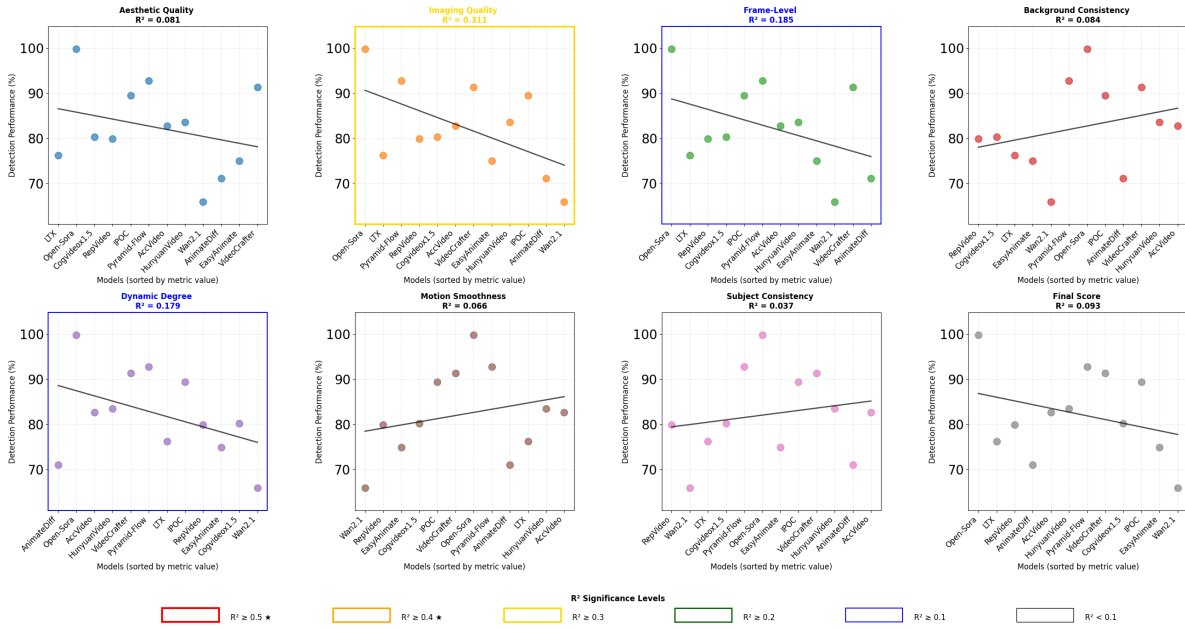
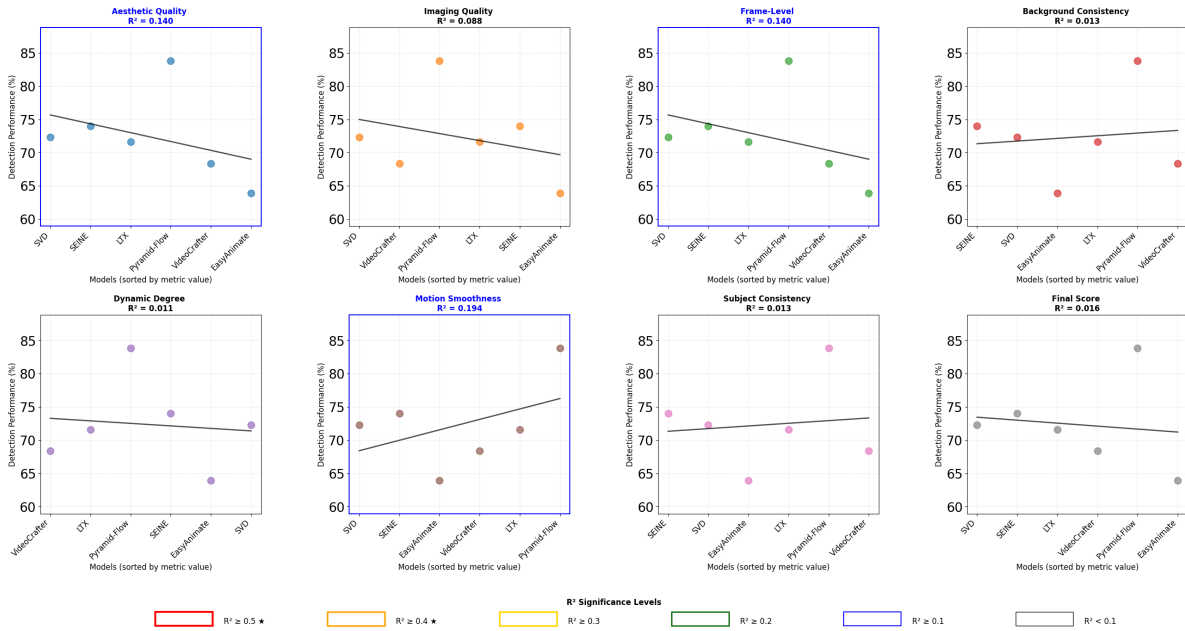


Figure 46. Correlations between the performance of TSM and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

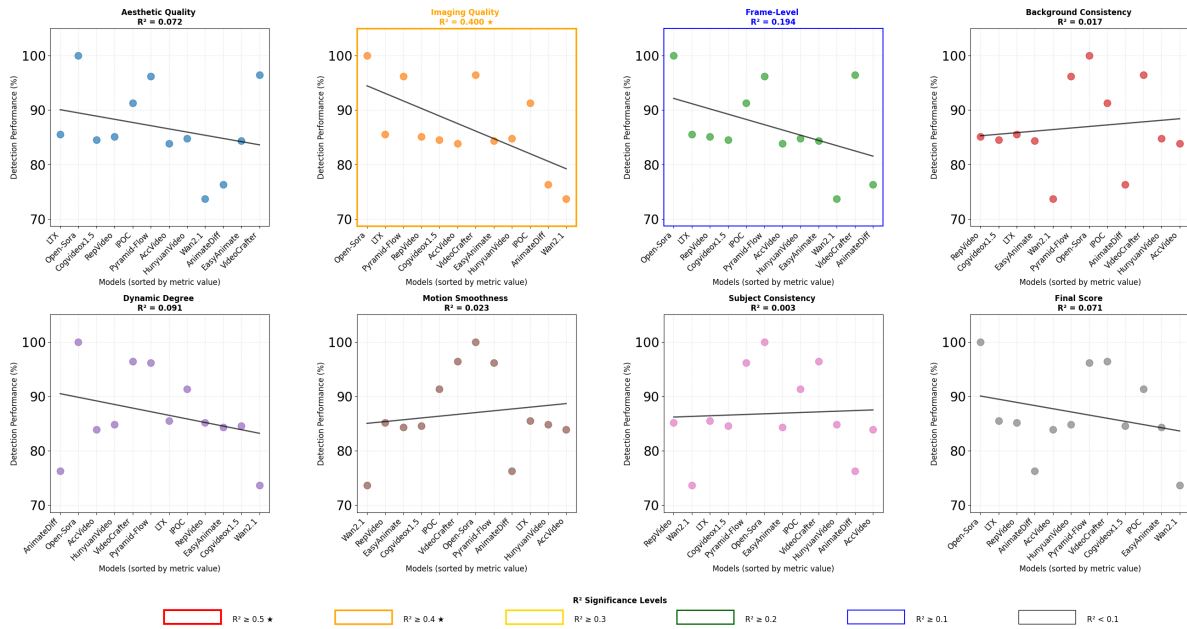


(a) Text-to-Video (T2V) Task.

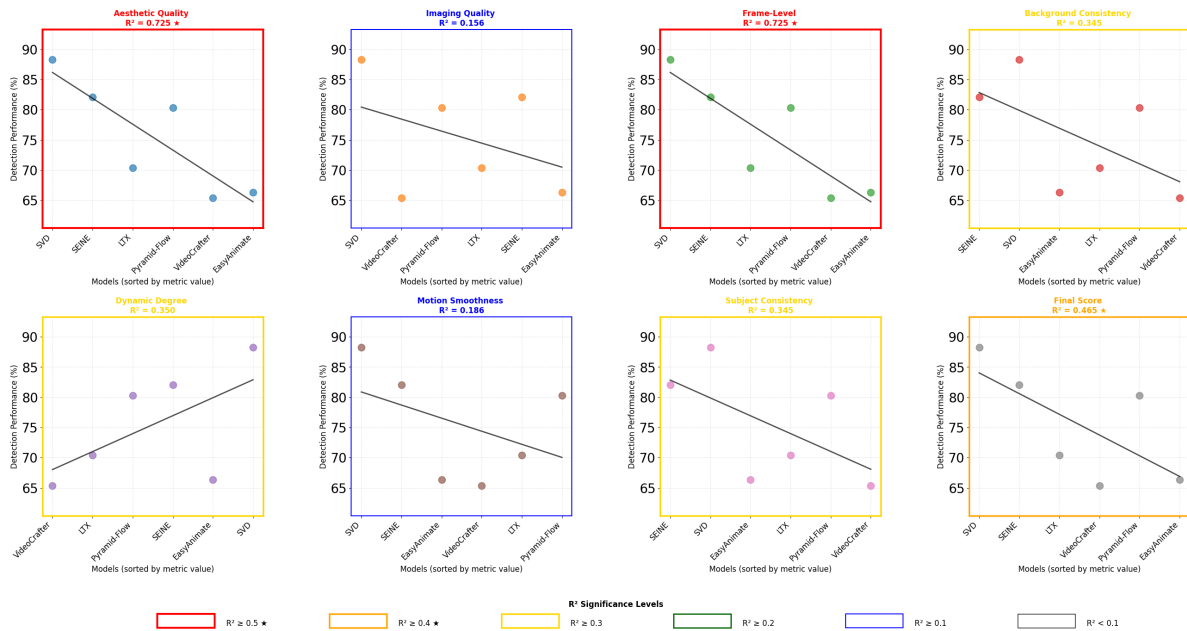


(b) Image-to-Video (I2V) Task.

Figure 47. Correlations between the performance of SlowFast and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).



(a) Text-to-Video (T2V) Task.



(b) Image-to-Video (I2V) Task.

Figure 48. Correlations between the performance of TimeSformer and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

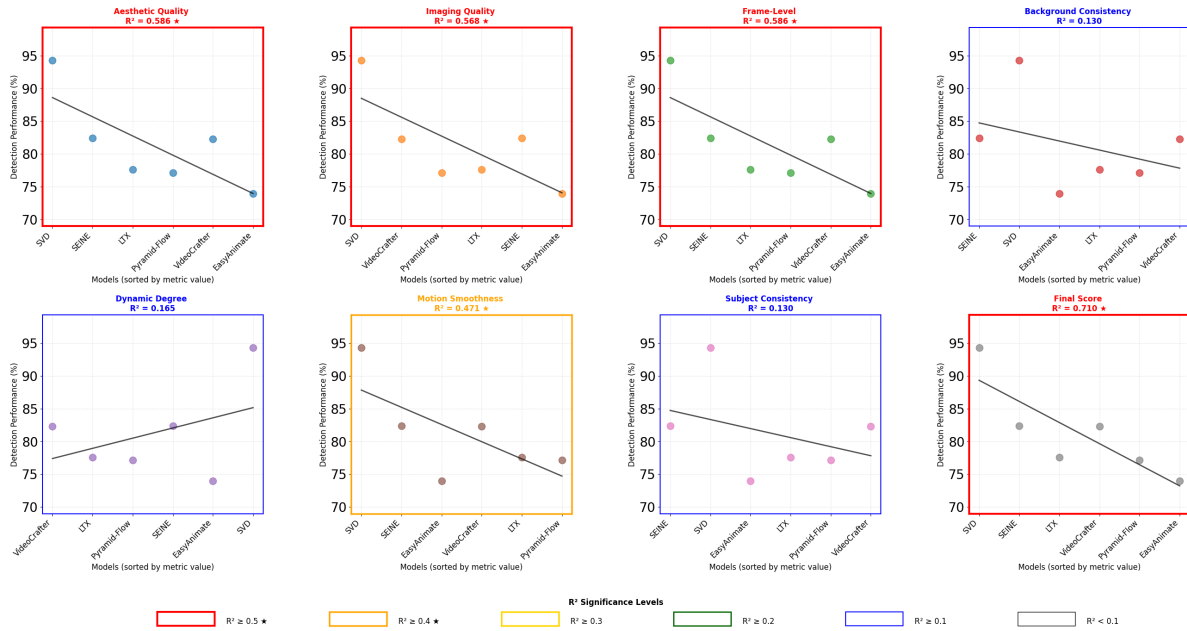
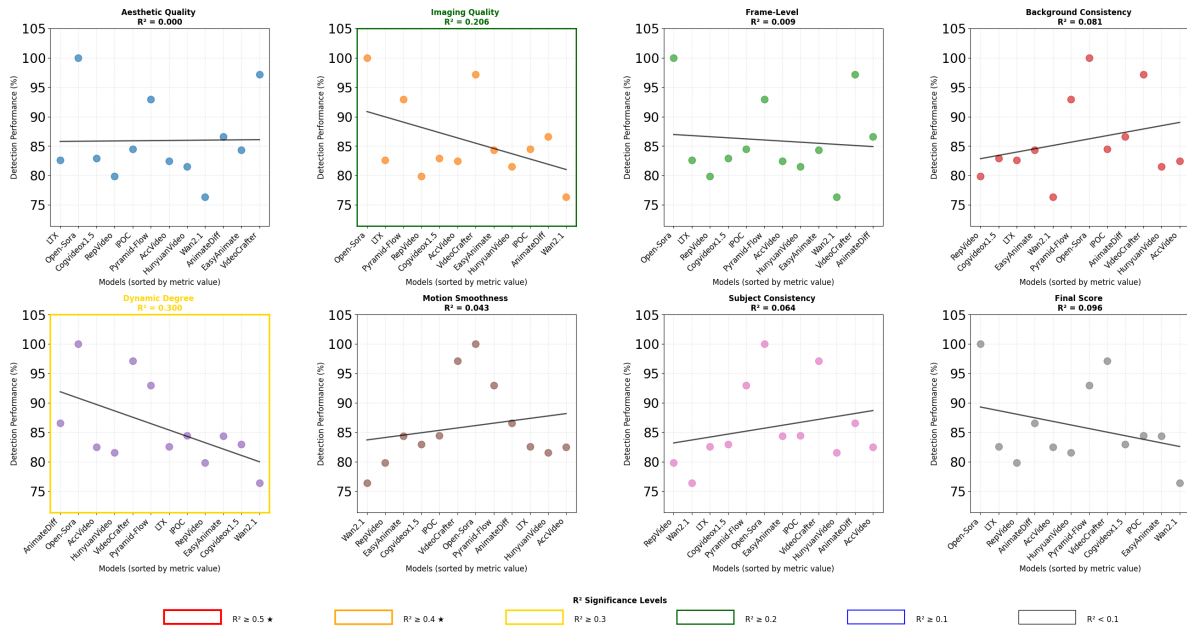


Figure 49. Correlations between the performance of UniFormerV2 and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

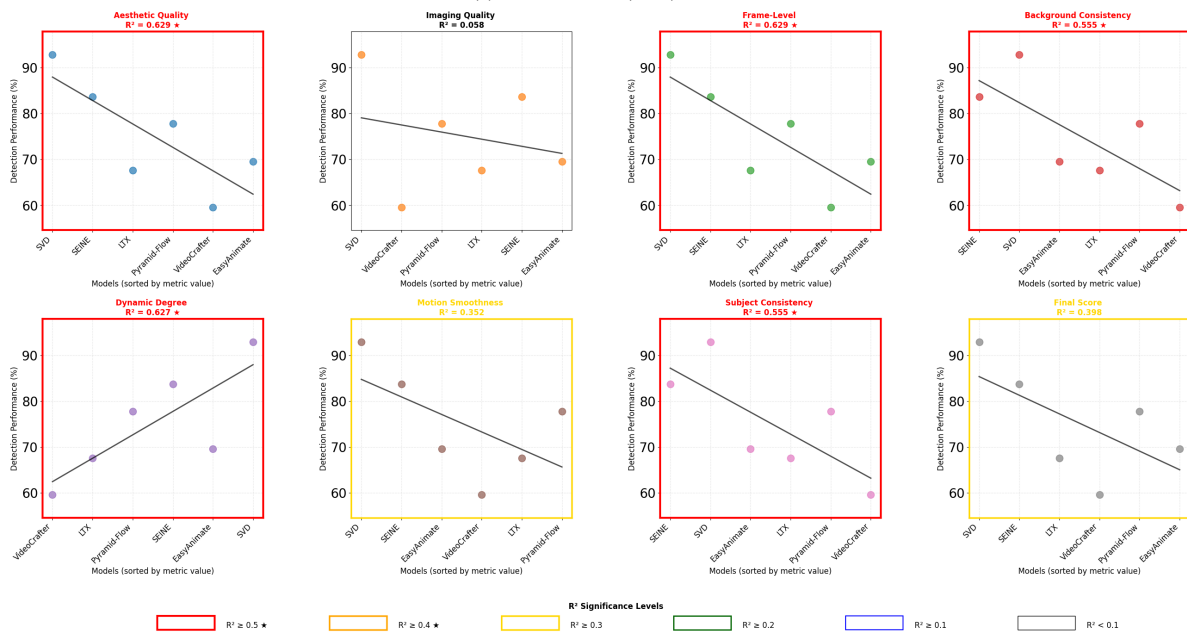
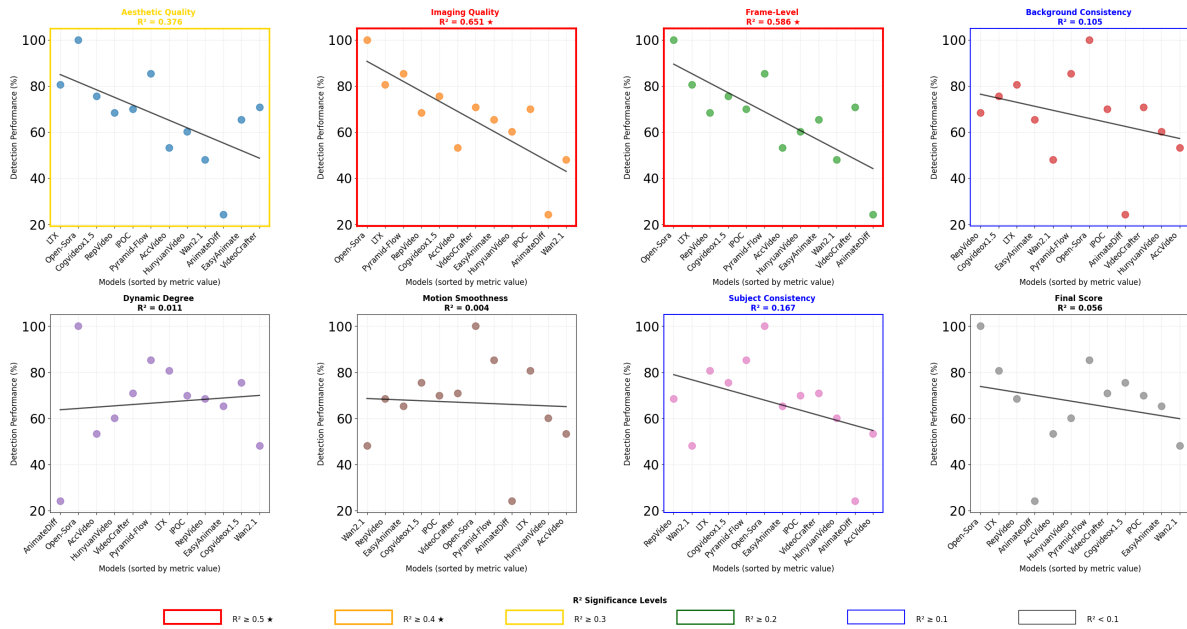


Figure 50. Correlations between the performance of NPR and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

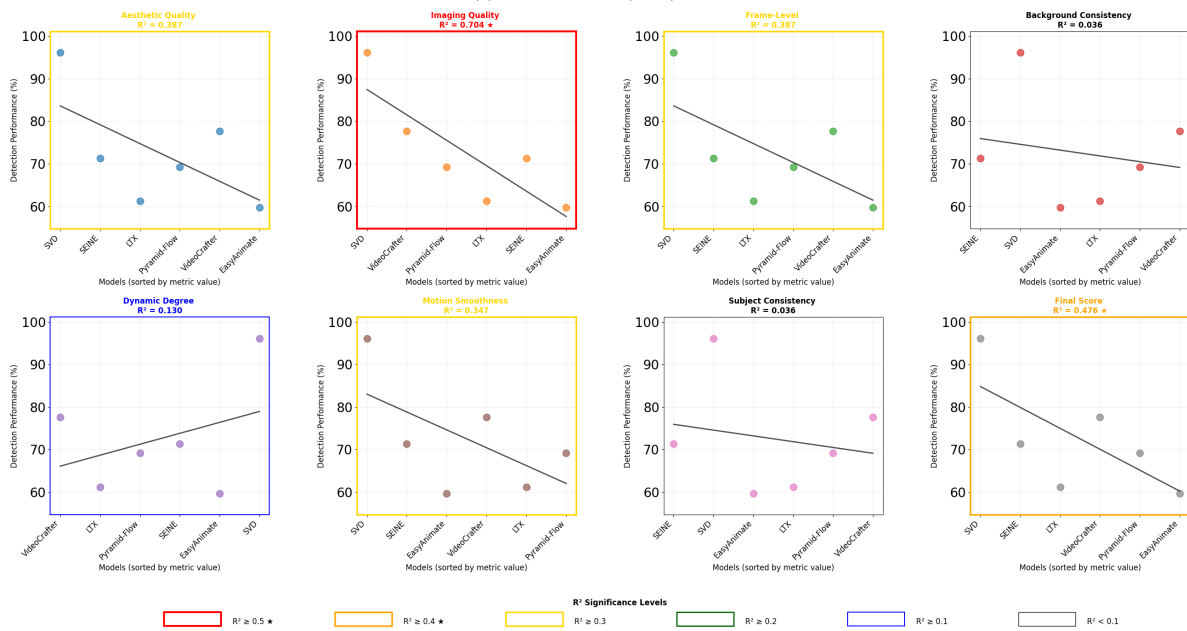
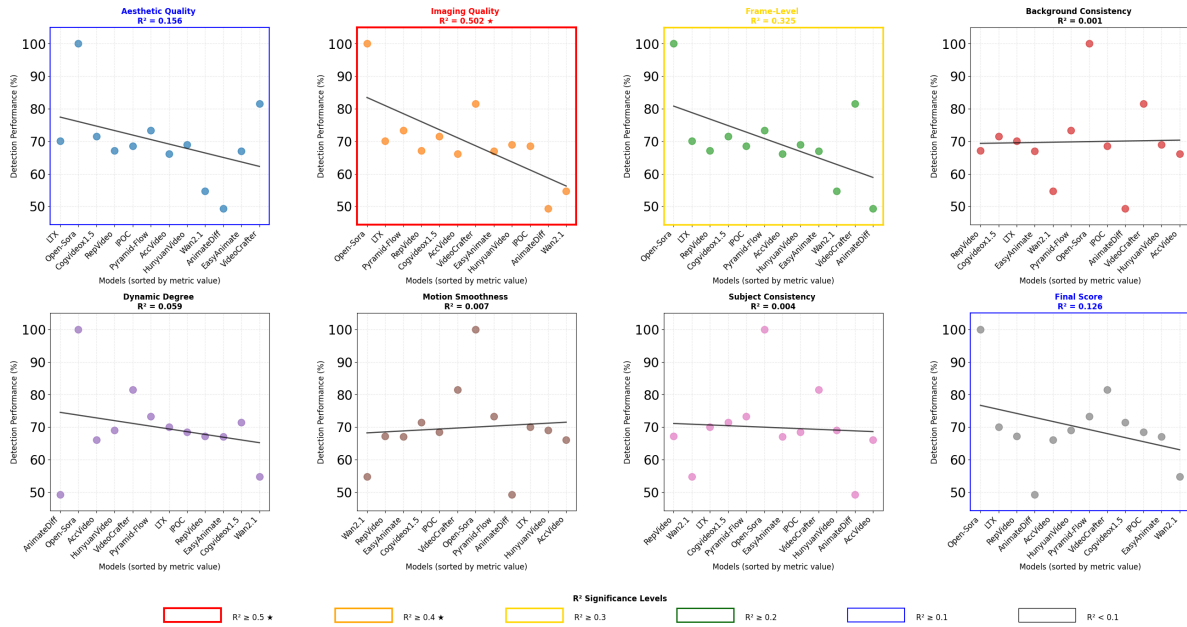


Figure 51. Correlations between the performance of FreDect and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

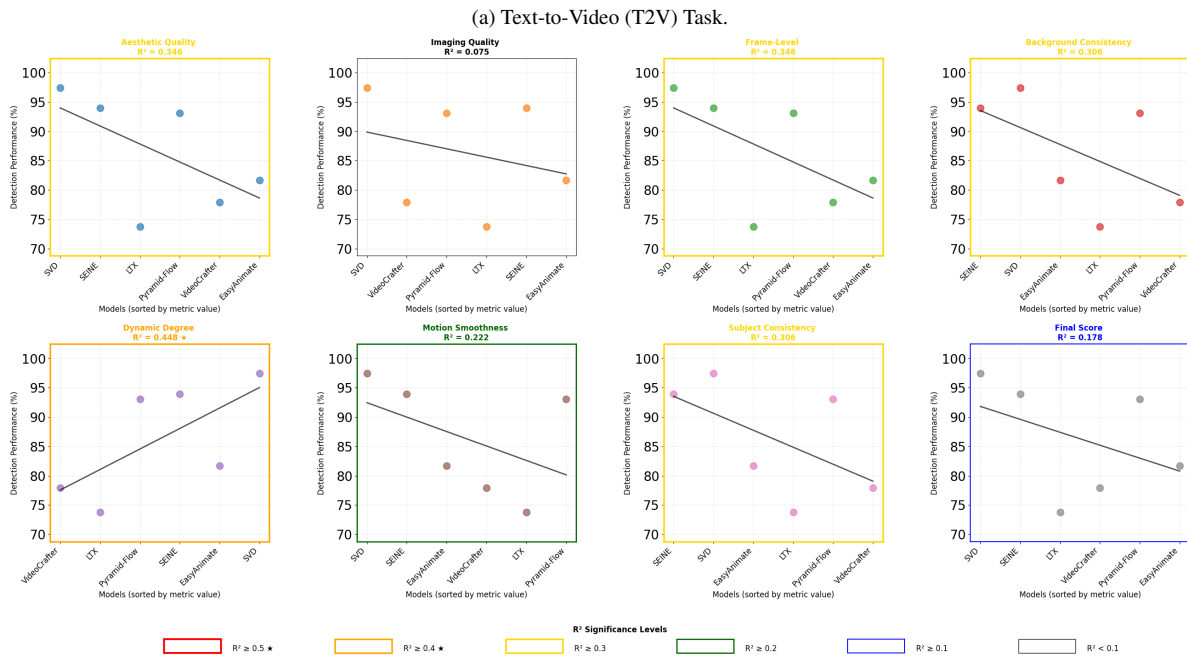
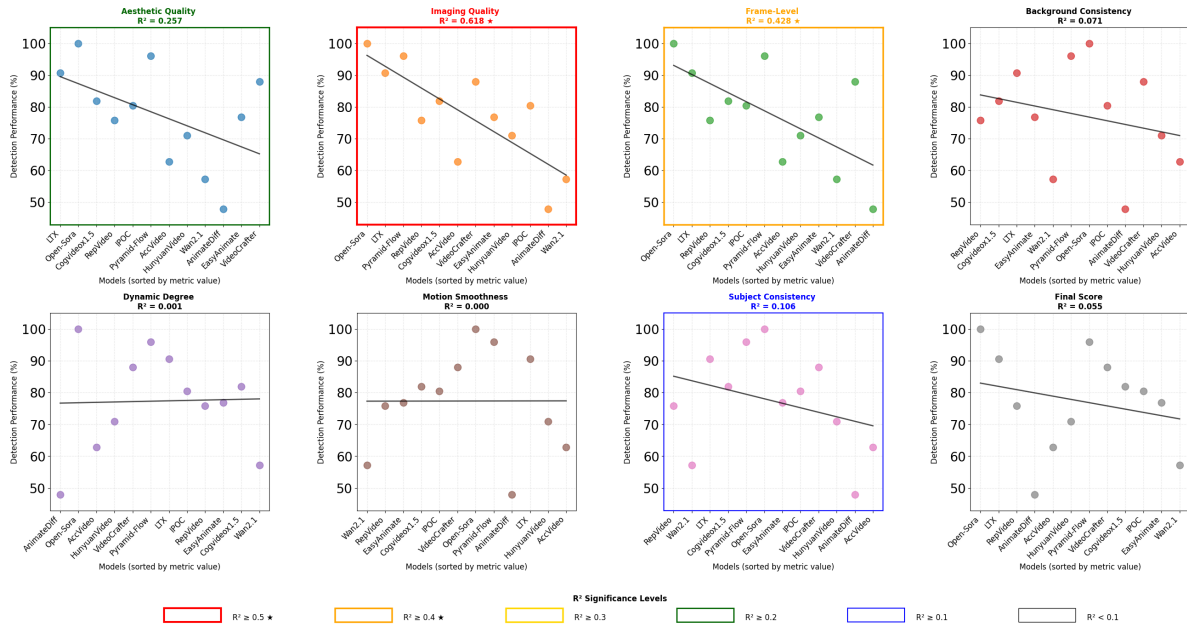


Figure 52. Correlations between the performance of Fusing and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

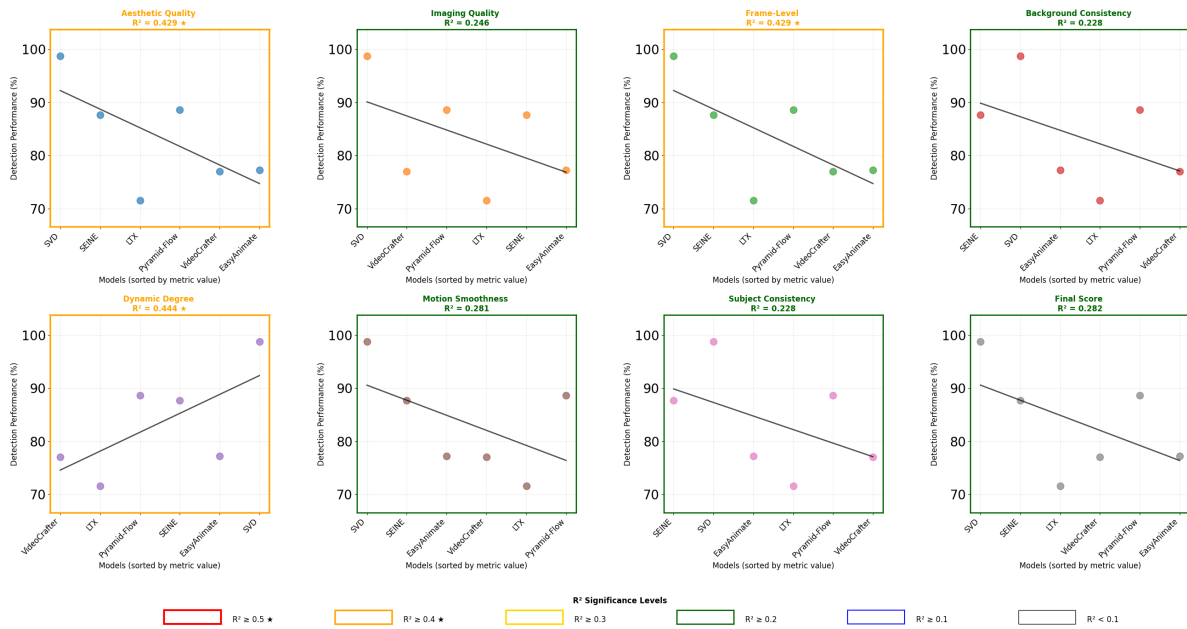
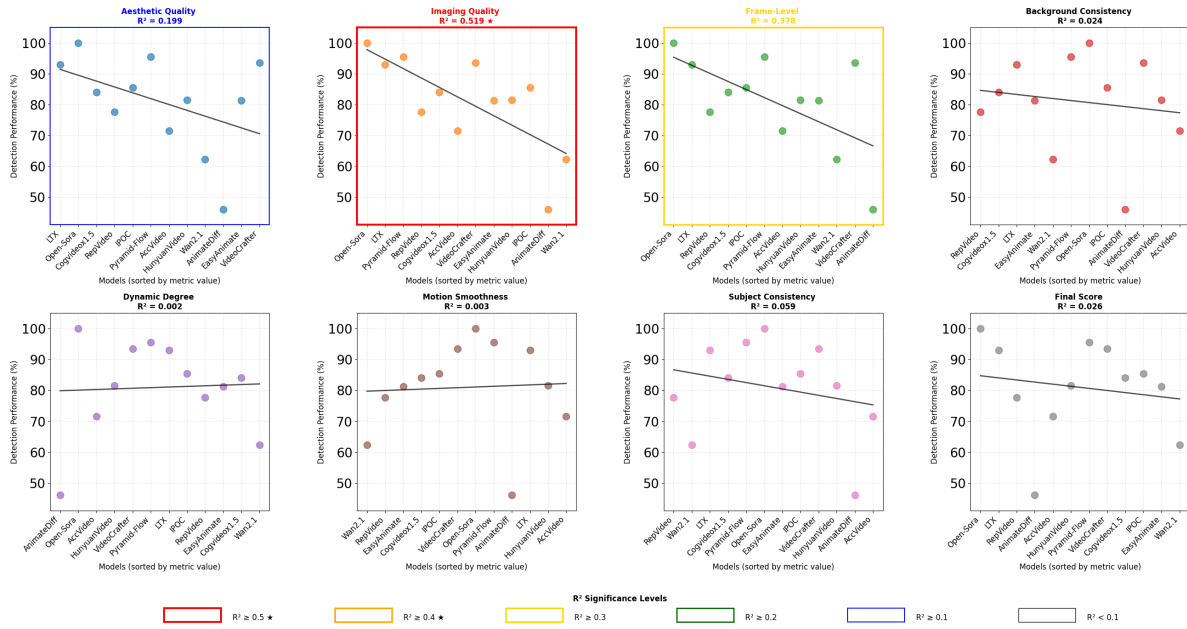


Figure 53. Correlations between the performance of Gram-Net and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

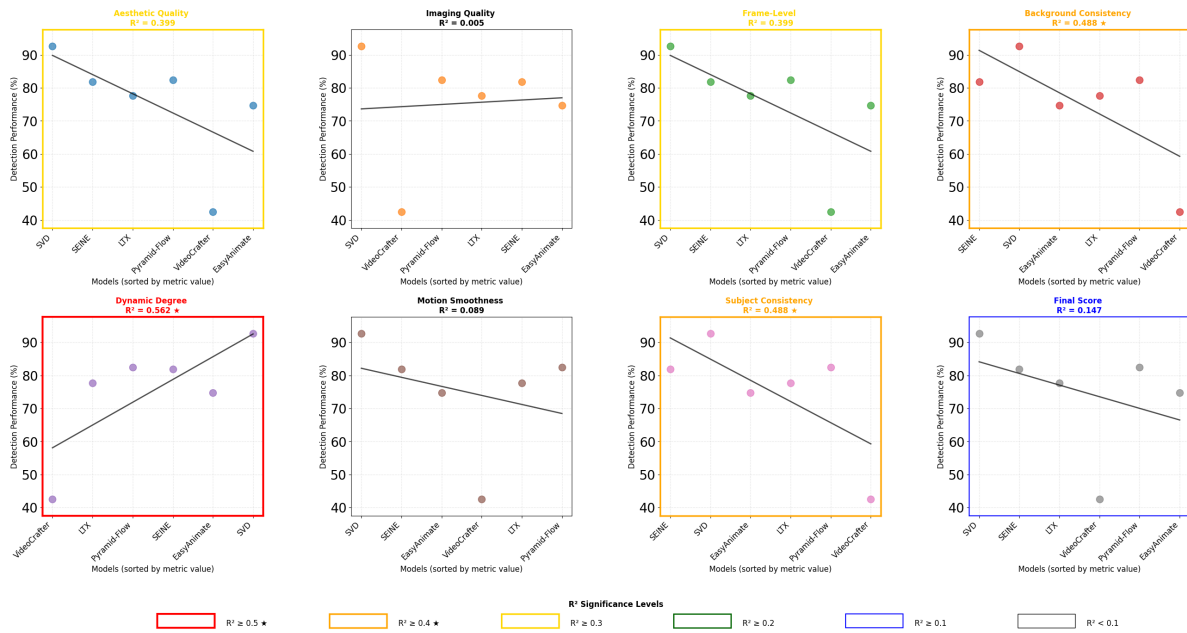
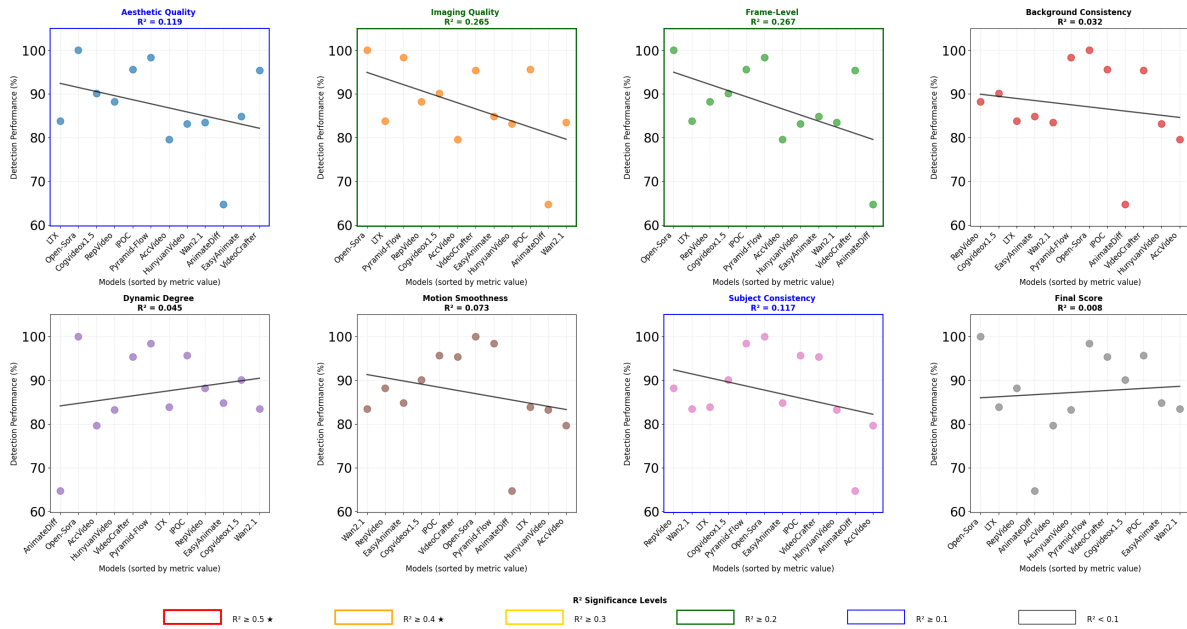


Figure 54. Correlations between the performance of D3 and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

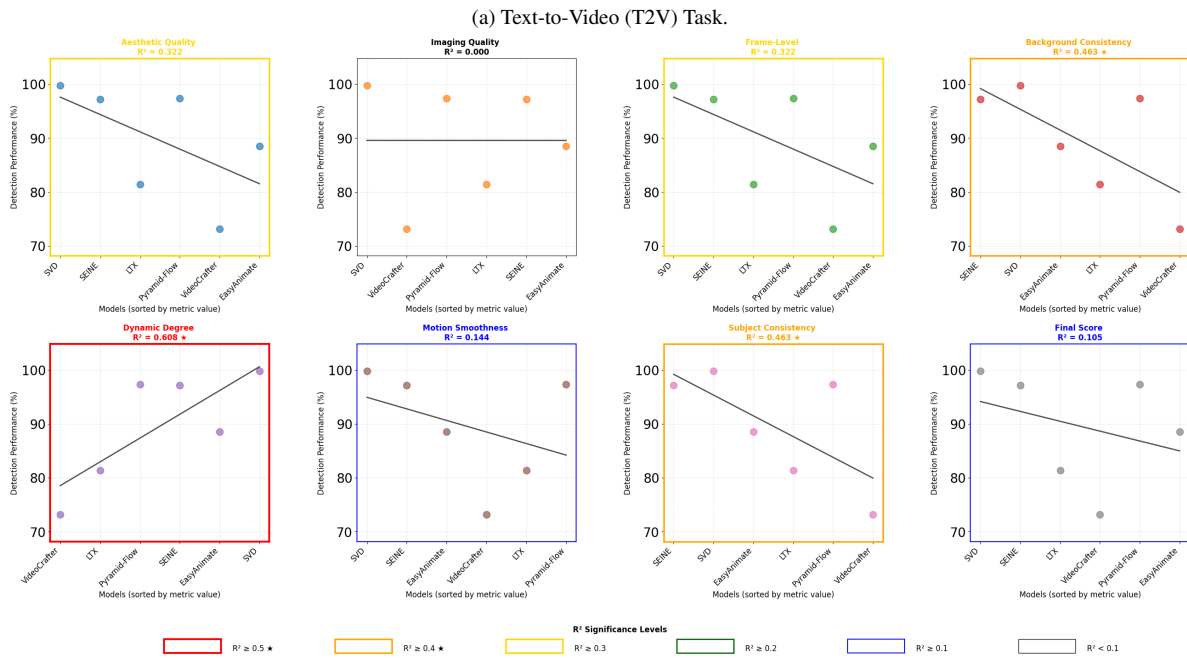
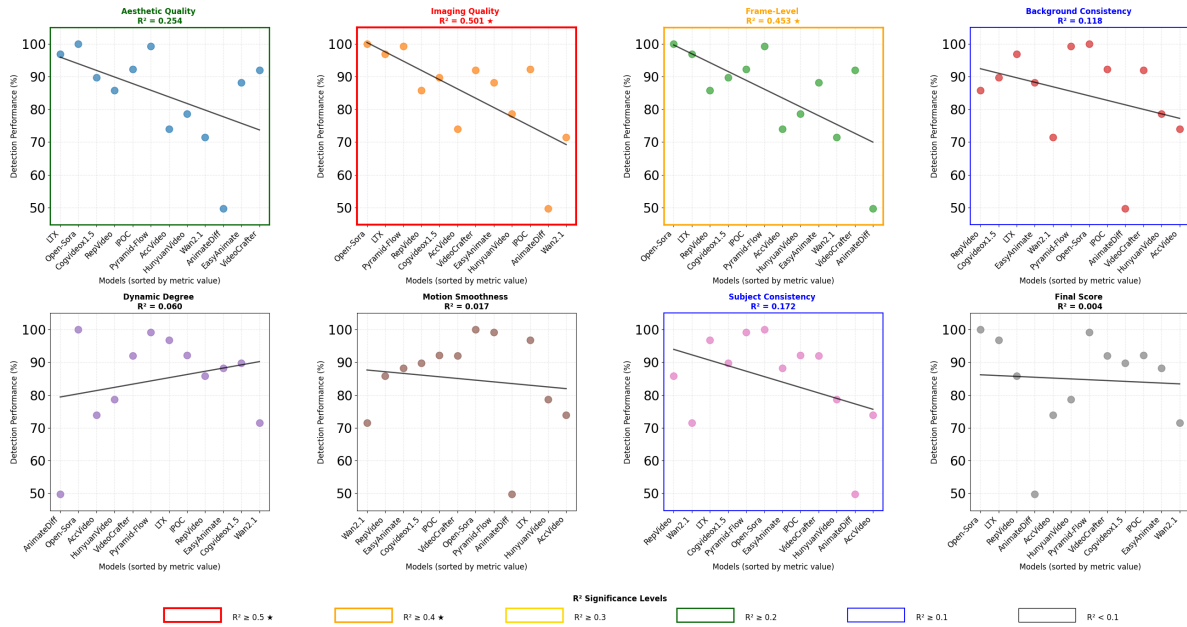


Figure 55. Correlations between the performance of CNNSpot and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

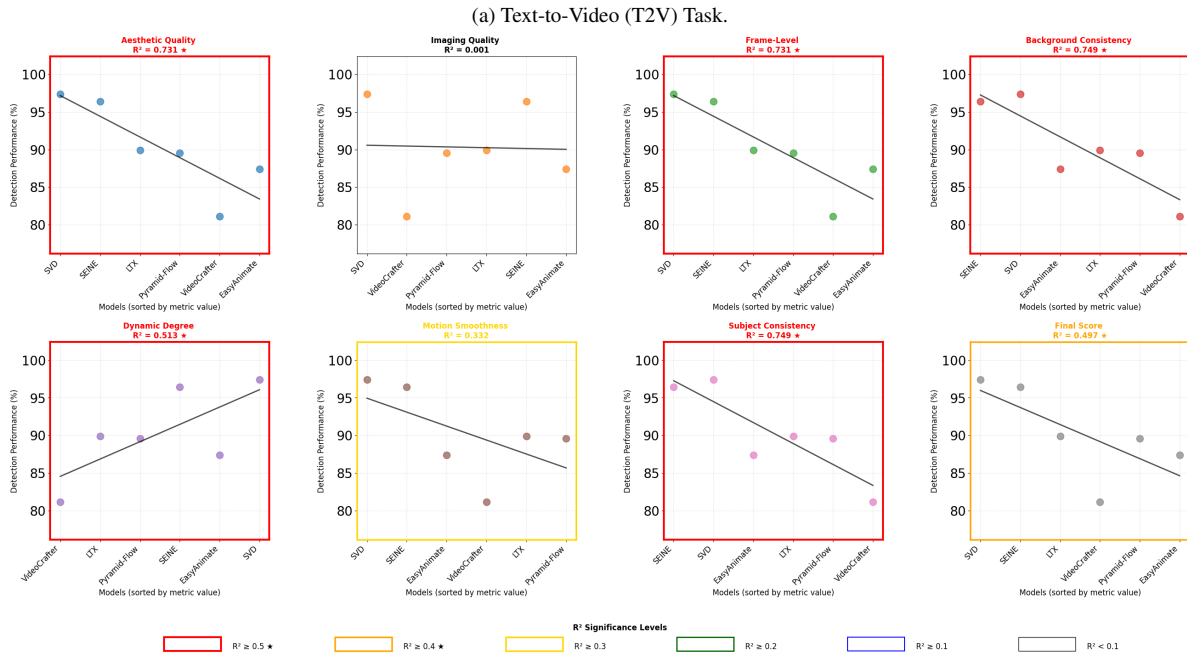
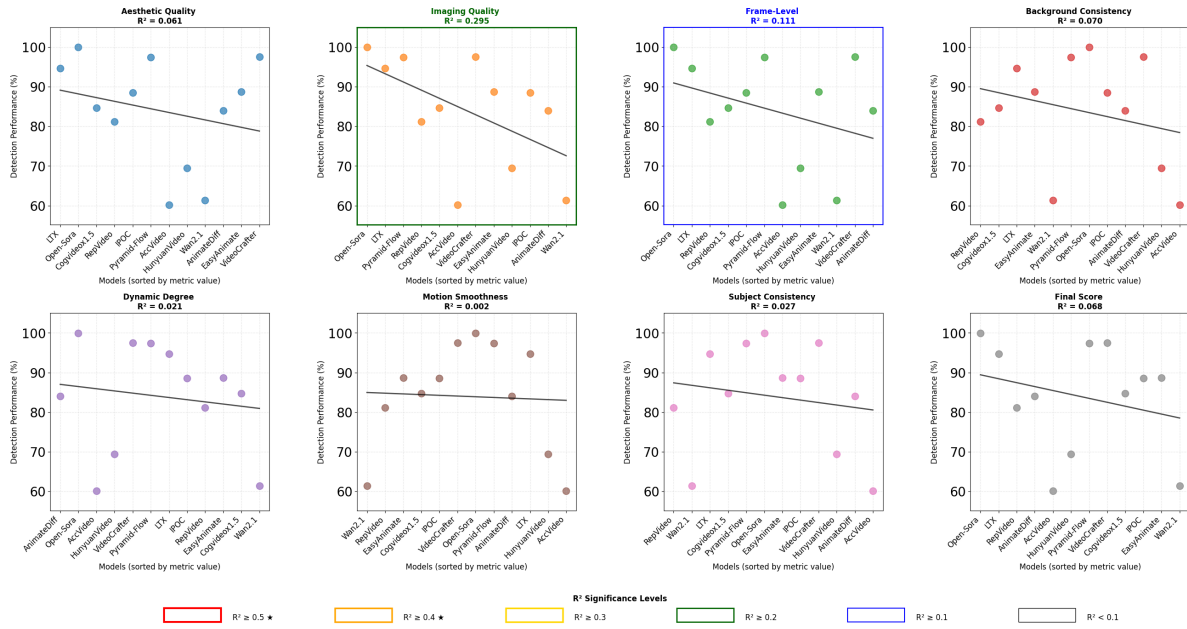


Figure 56. Correlations between the performance of UnivFD and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

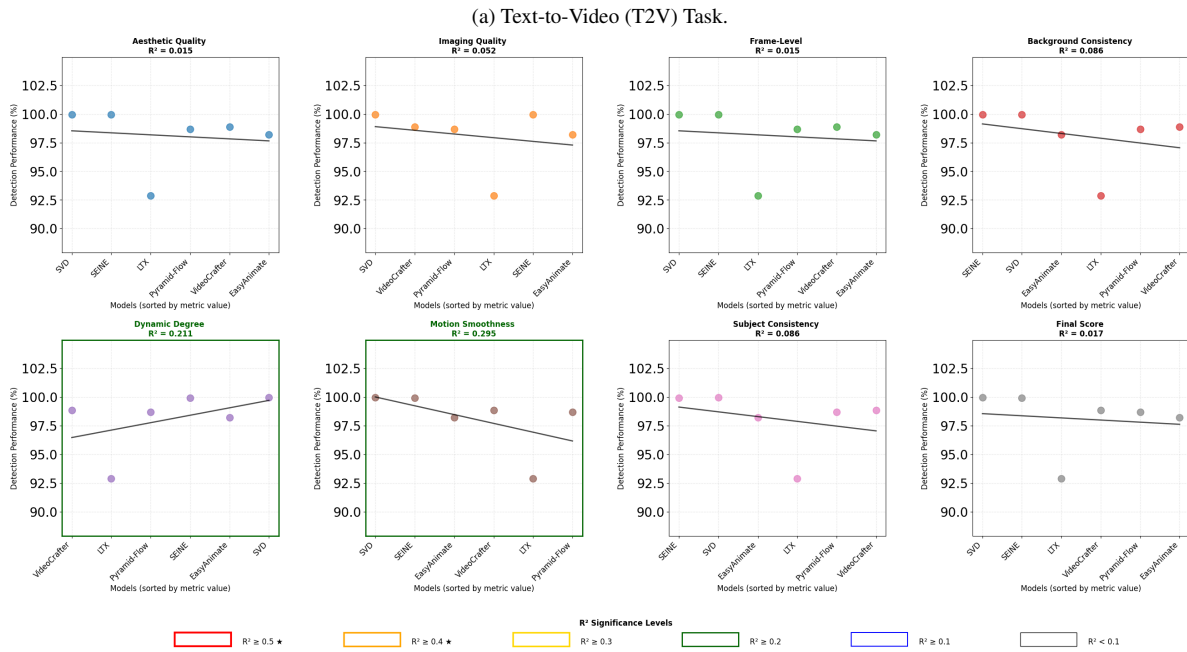
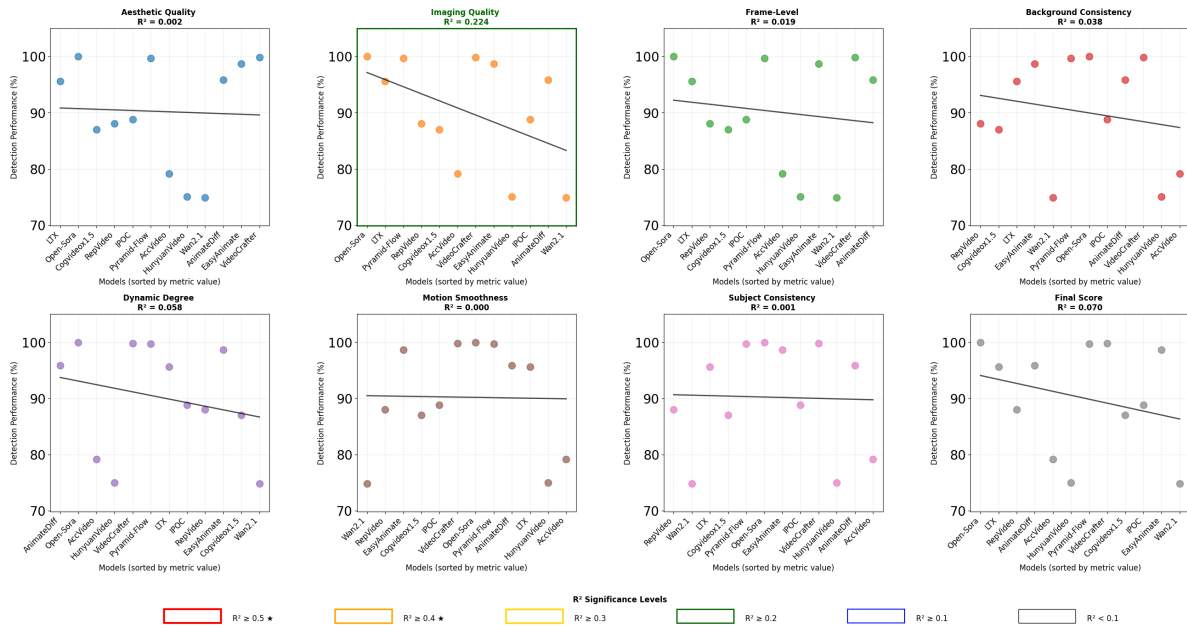


Figure 57. Correlations between the performance of ForgeLens1 and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

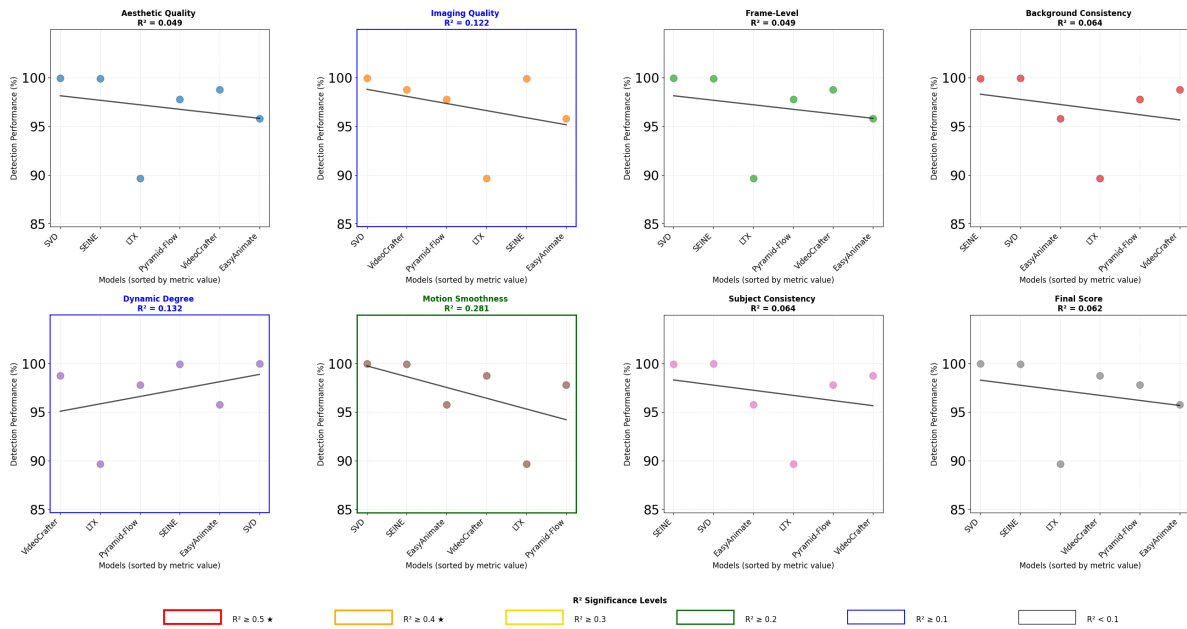
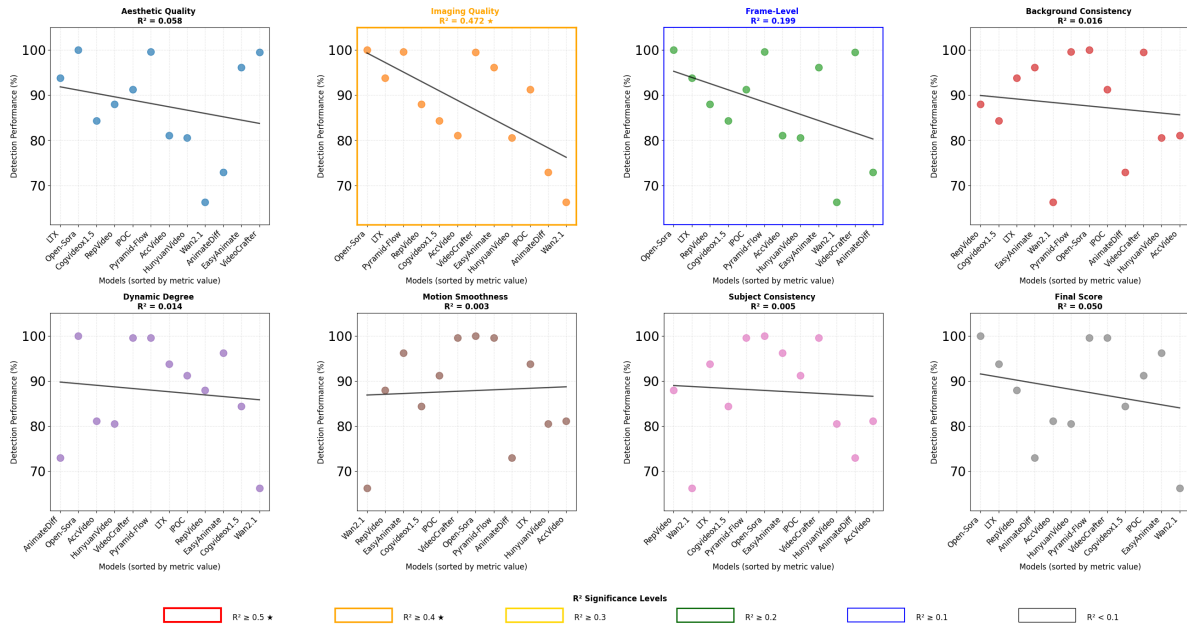


Figure 58. Correlations between the performance of ForgeLens3 and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

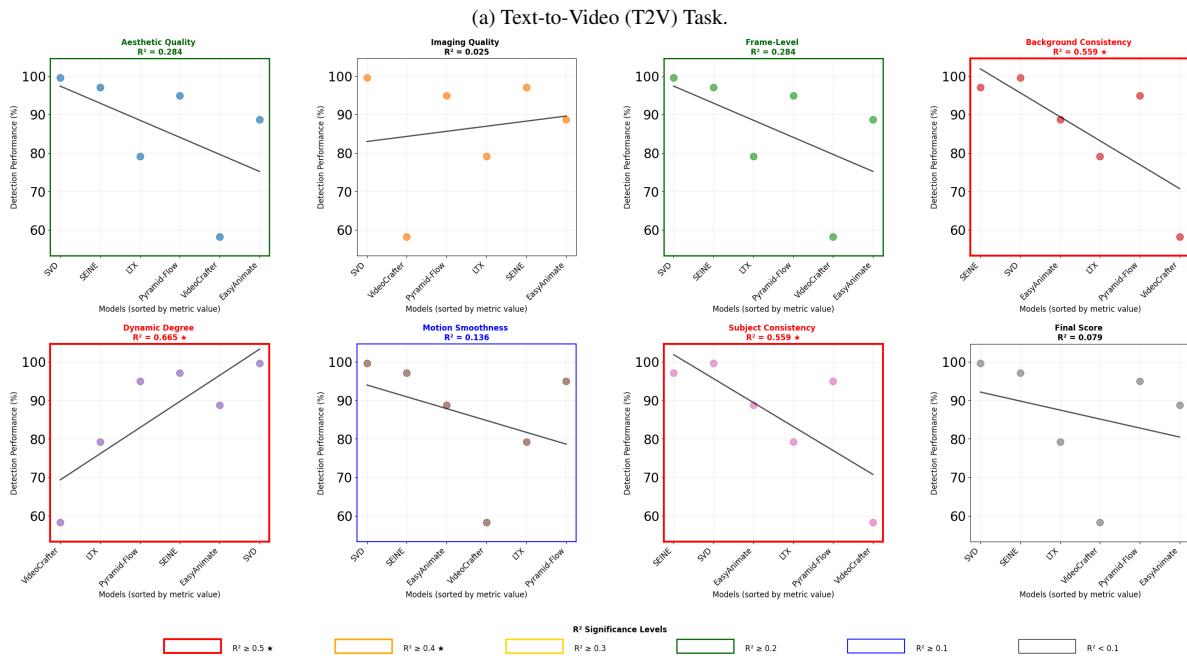
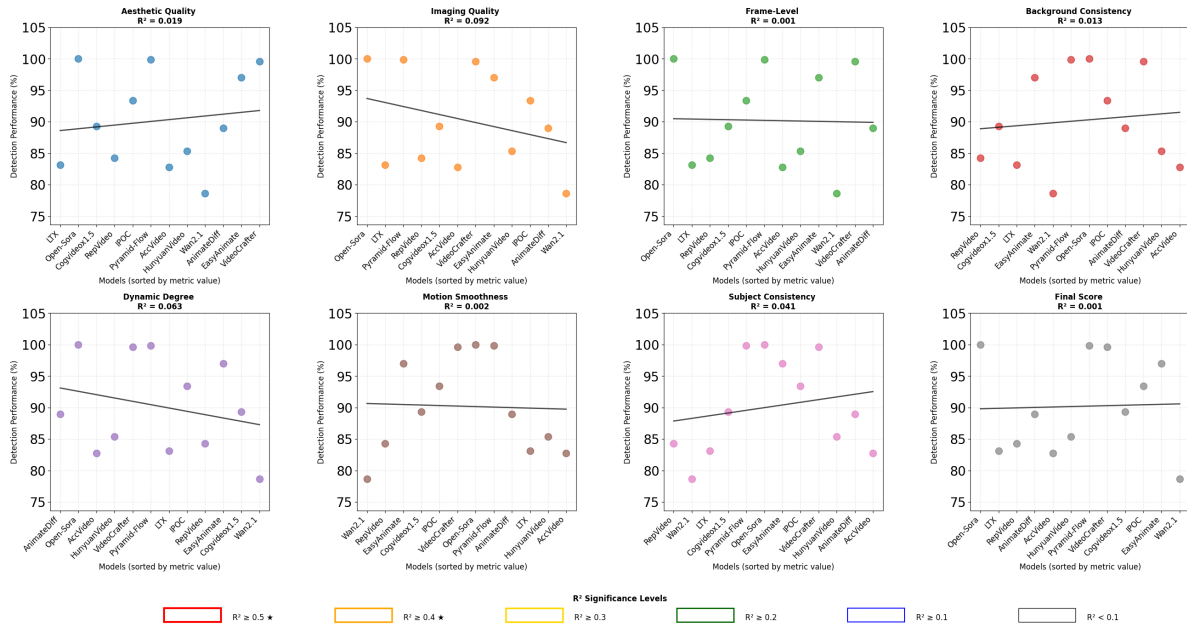


Figure 59. Correlations between the performance of Effort and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

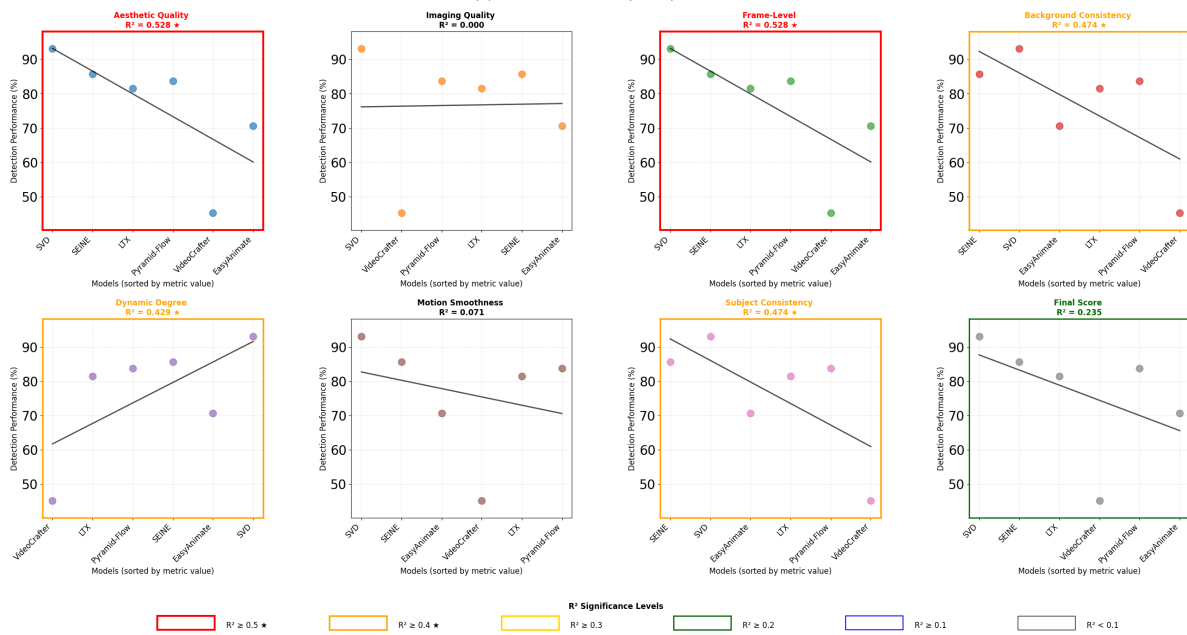
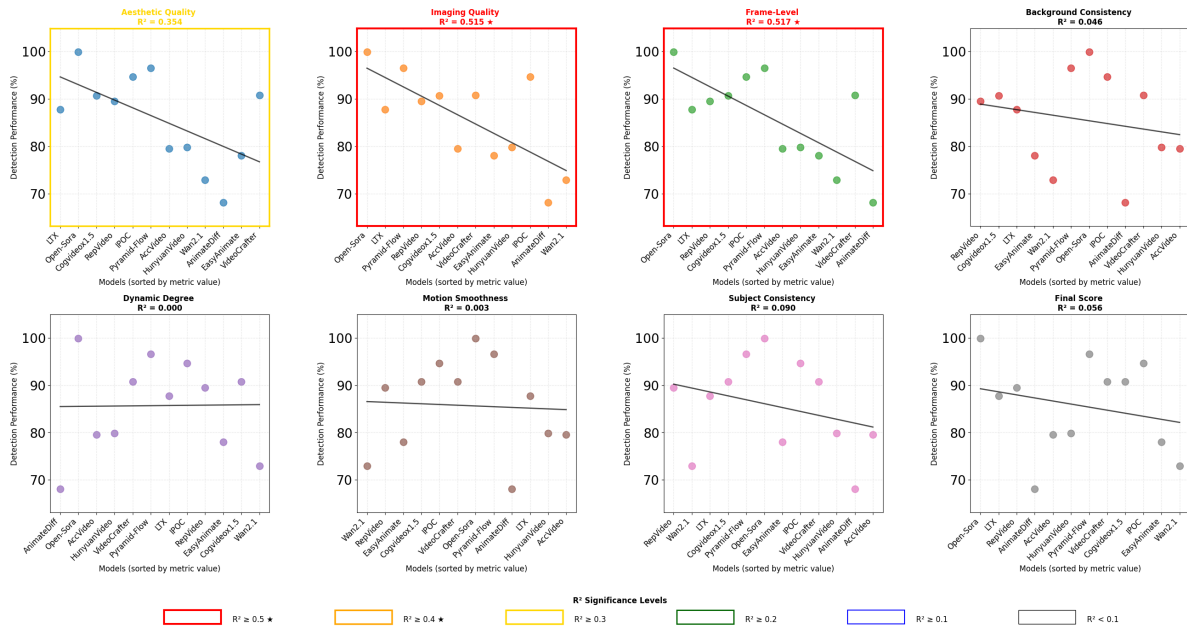
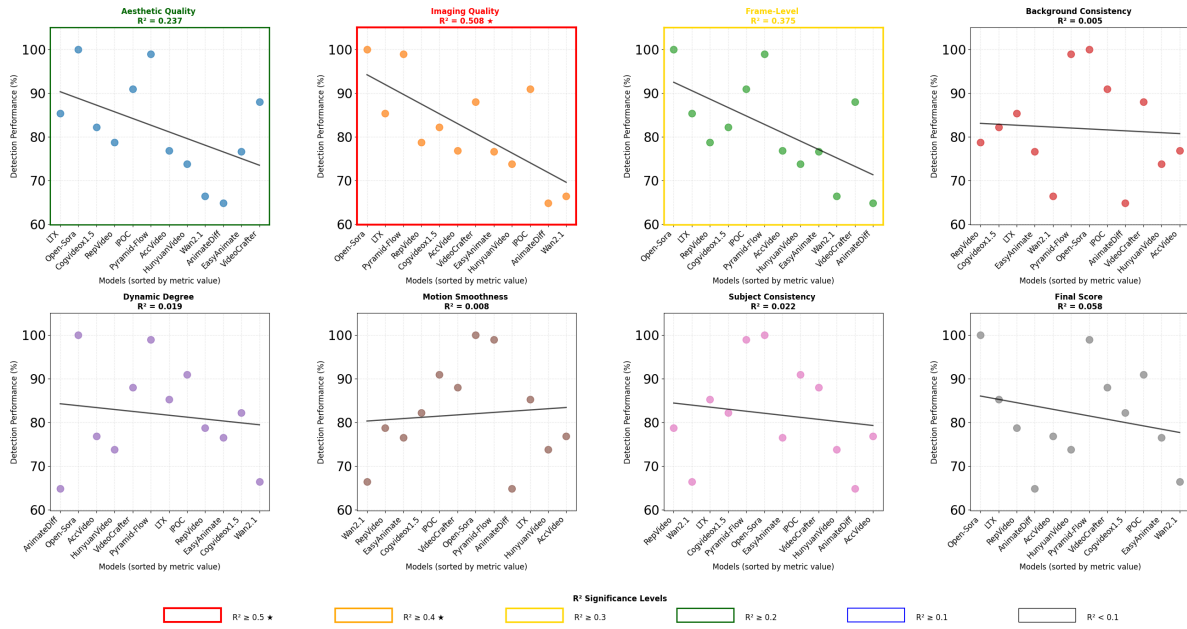
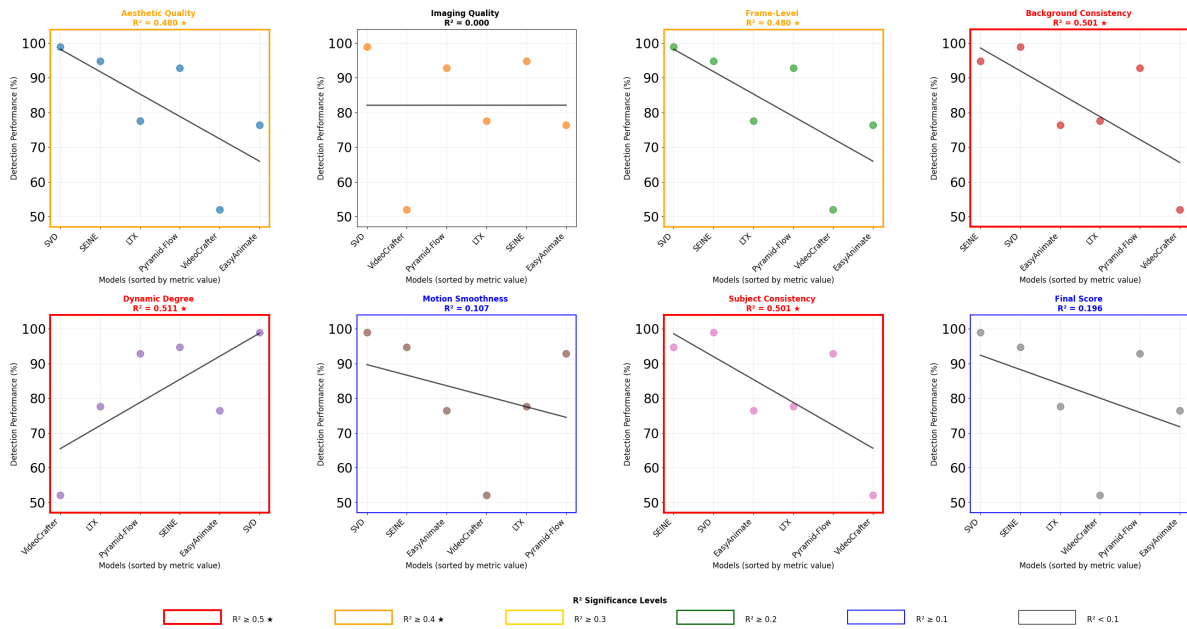


Figure 60. Correlations between the performance of DeCoF and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

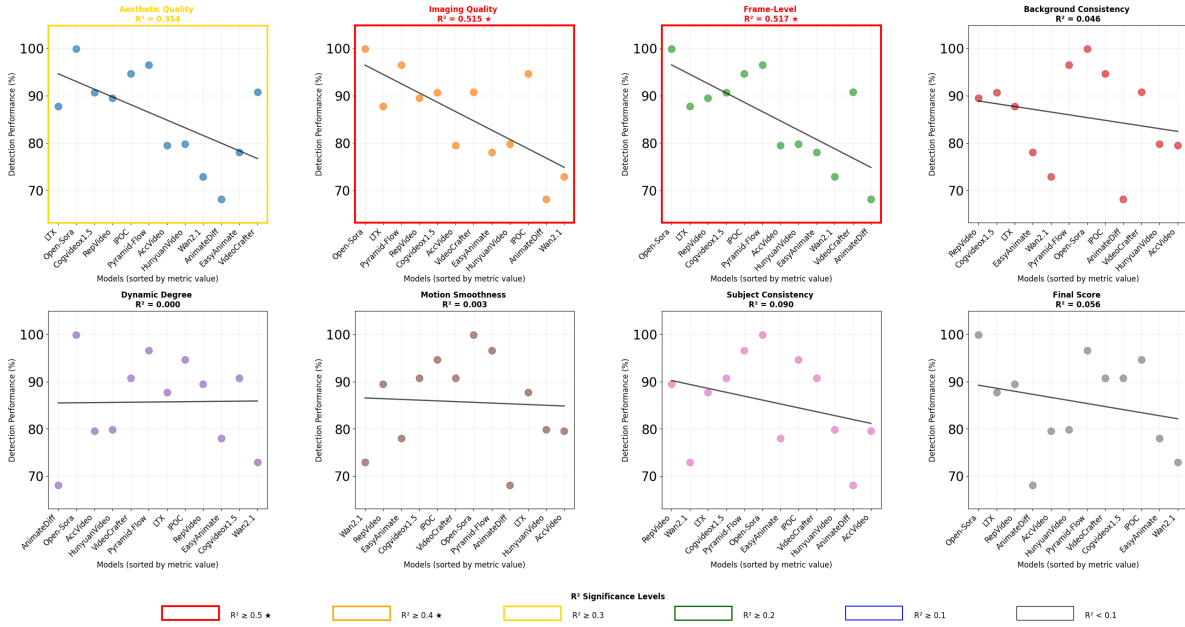


(a) Text-to-Video (T2V) Task.

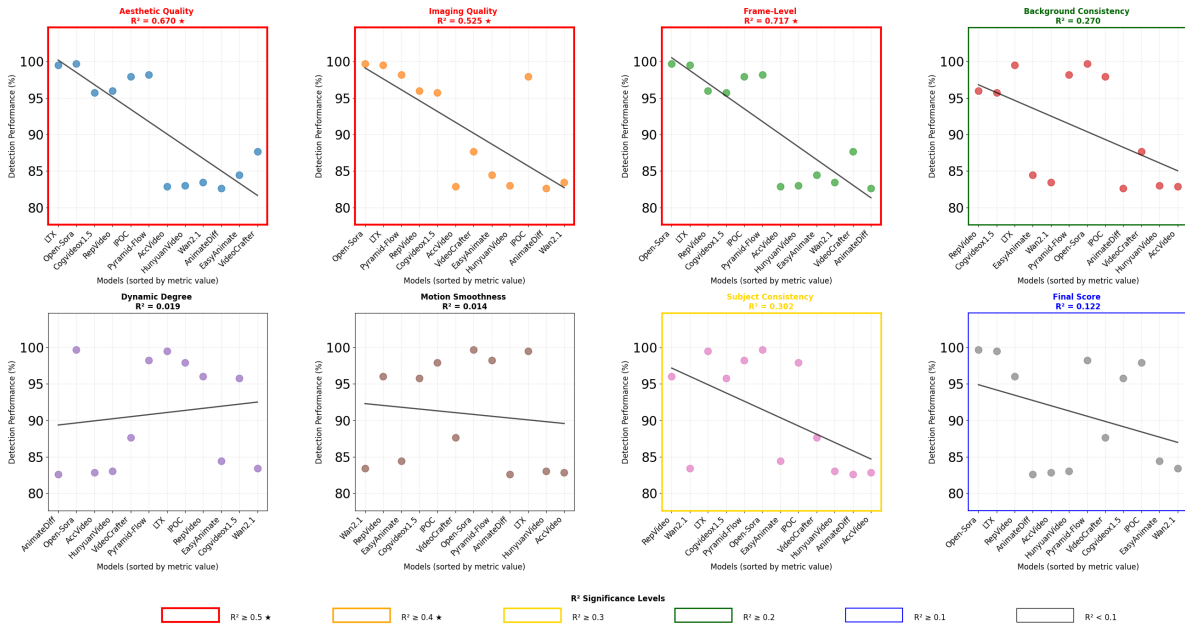


(b) Image-to-Video (I2V) Task.

Figure 61. Correlations between the performance of DeMamba and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

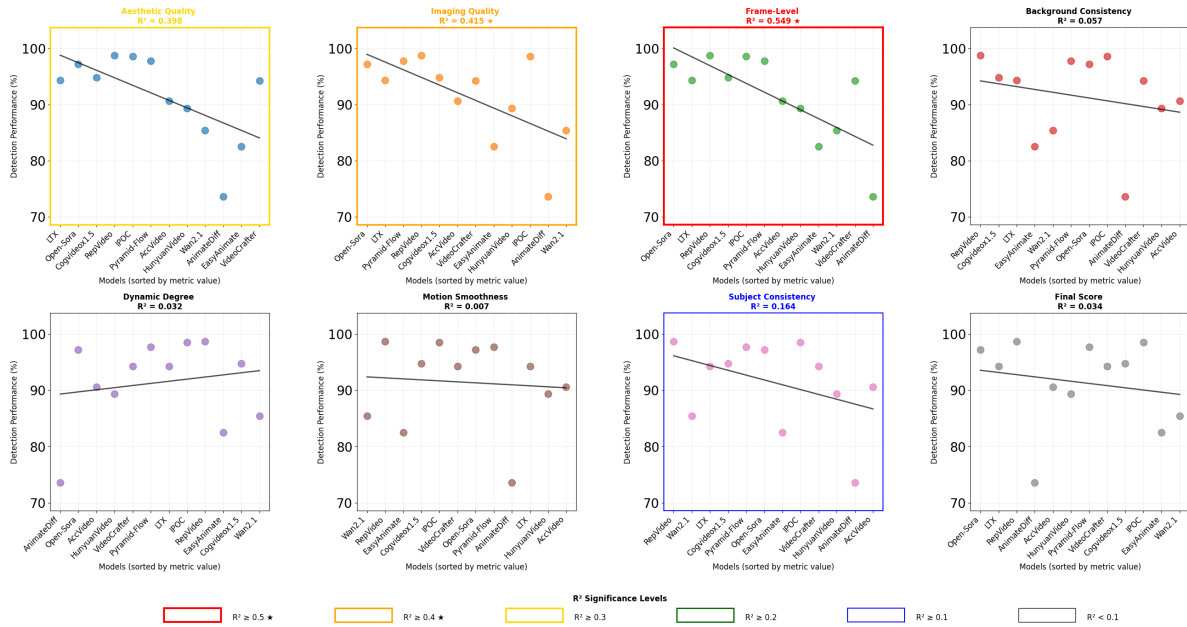


(a) Trained on Open-Sora.

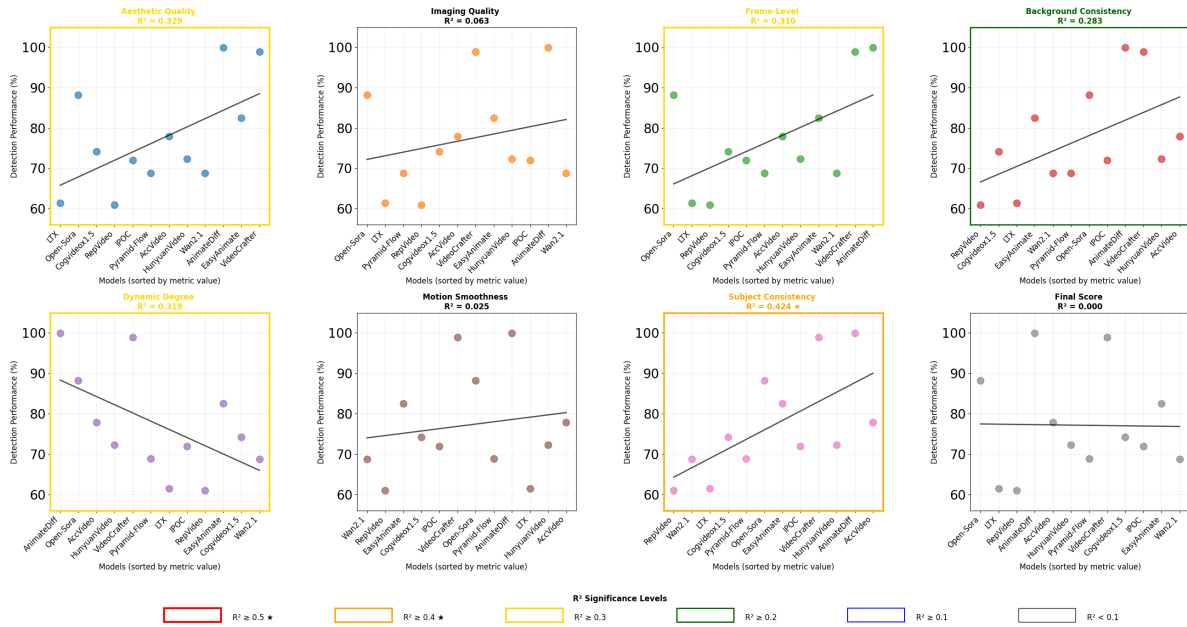


(b) Trained on LTX.

Figure 62. Correlations between the performance of DeCoF and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

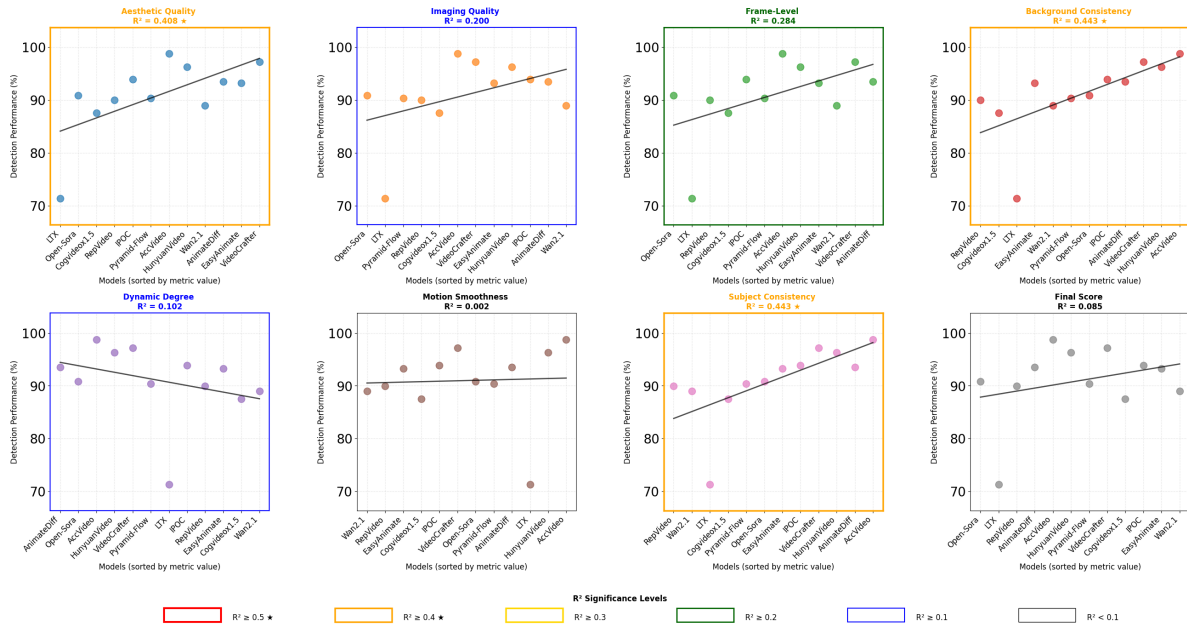


(a) Trained on RepVideo.

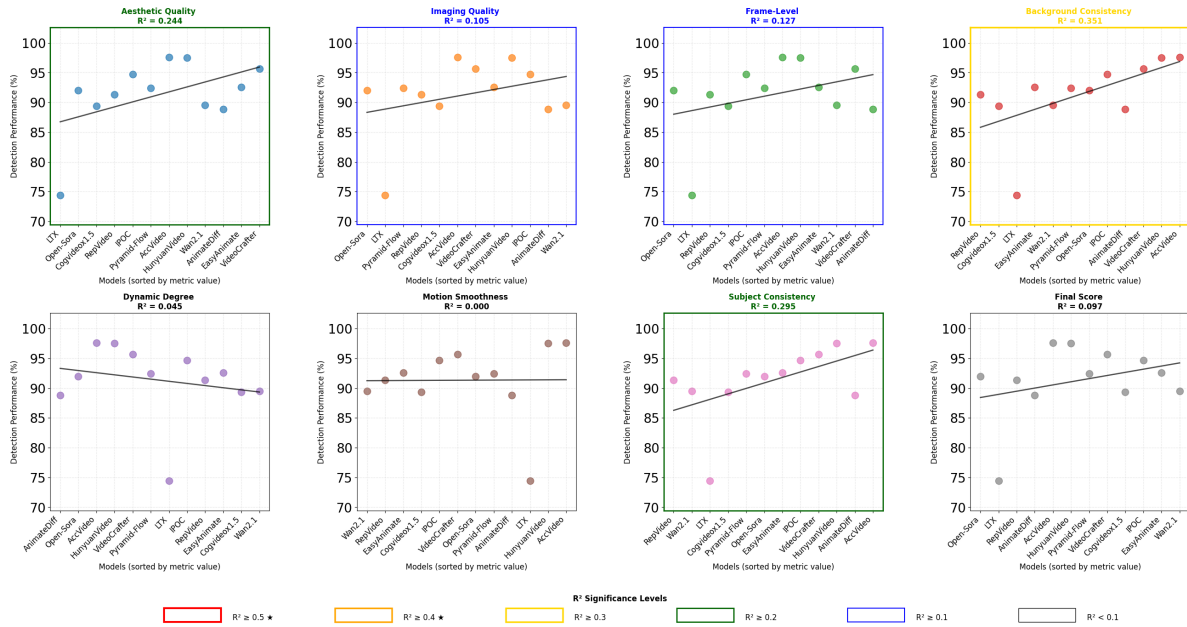


(b) Trained on AnimateDiff.

Figure 63. Correlations between the performance of DeCoF and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

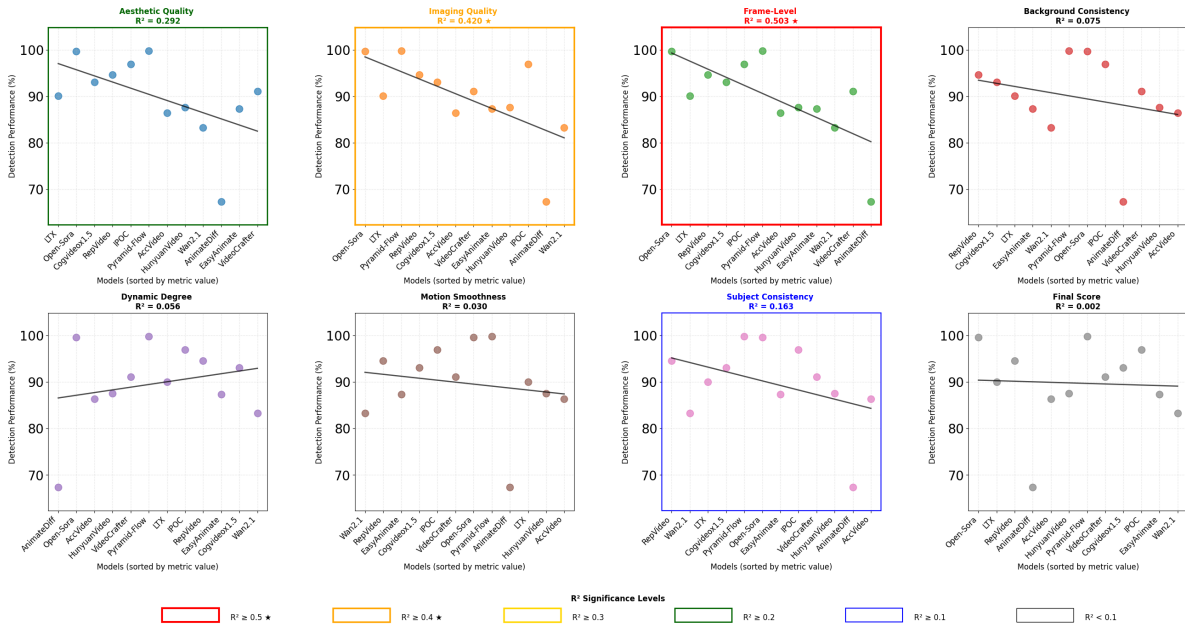


(a) Trained on AccVideo.

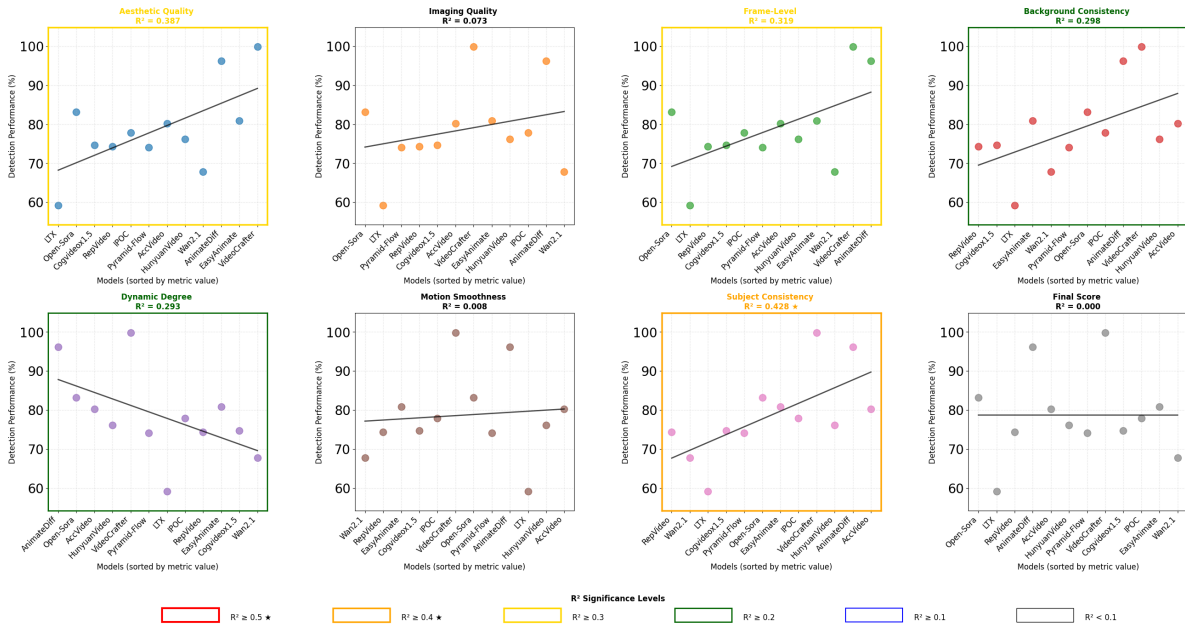


(b) Trained on HunyuanVideo.

Figure 64. Correlations between the performance of DeCoF and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

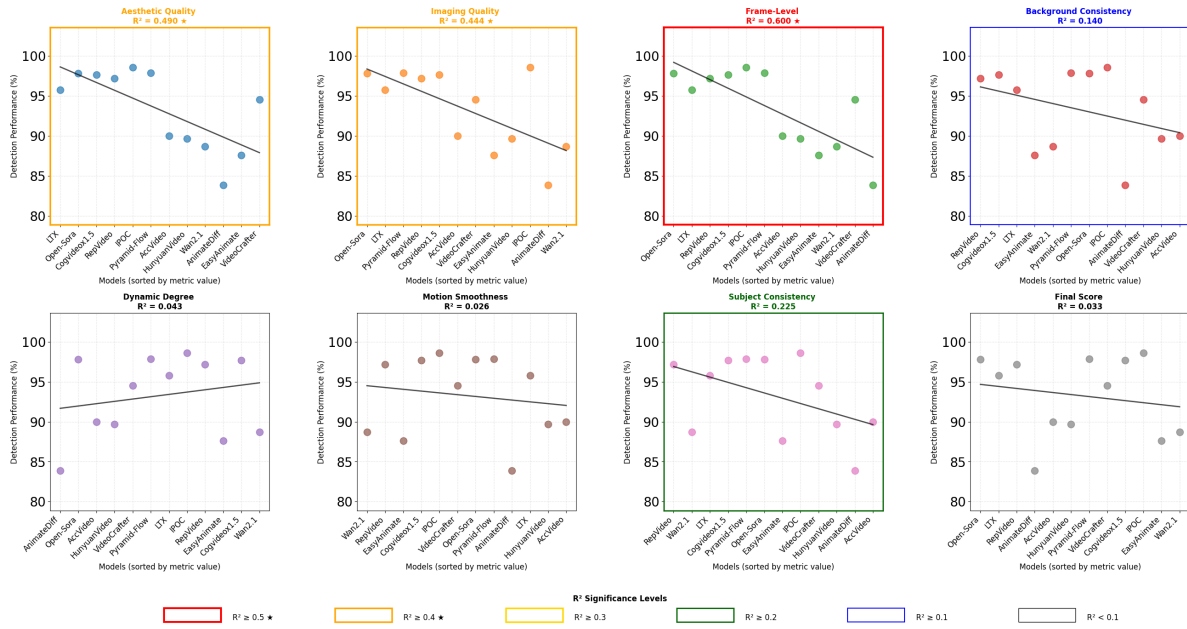


(a) Trained on Pyramid-Flow.

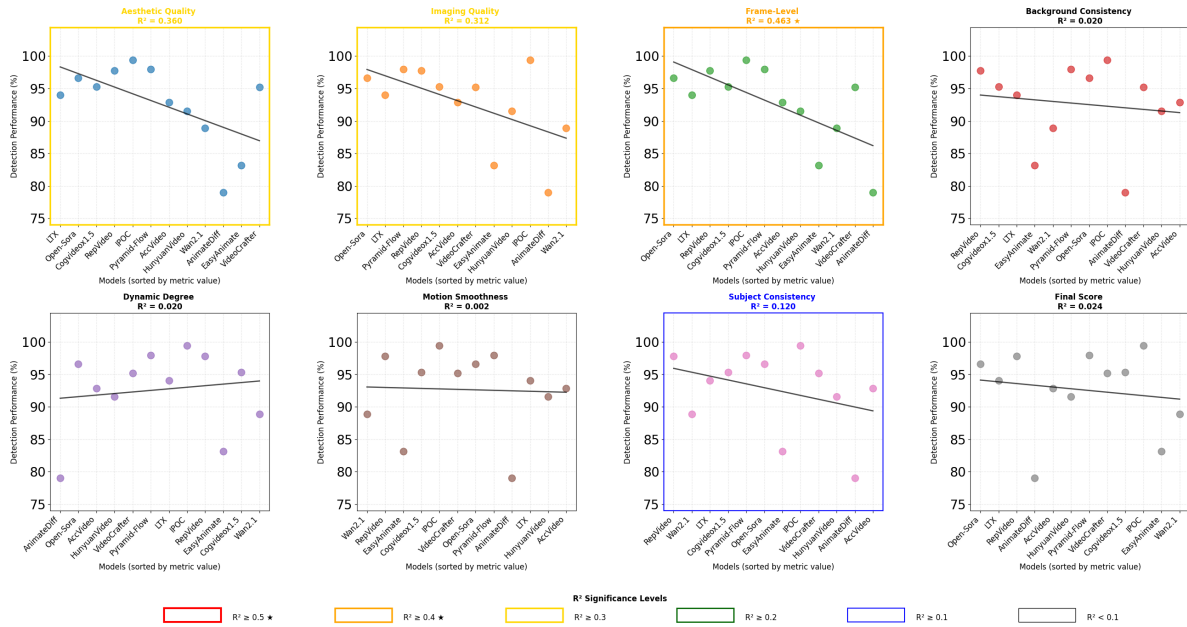


(b) Trained on VideoCrafter.

Figure 65. Correlations between the performance of DeCoF and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

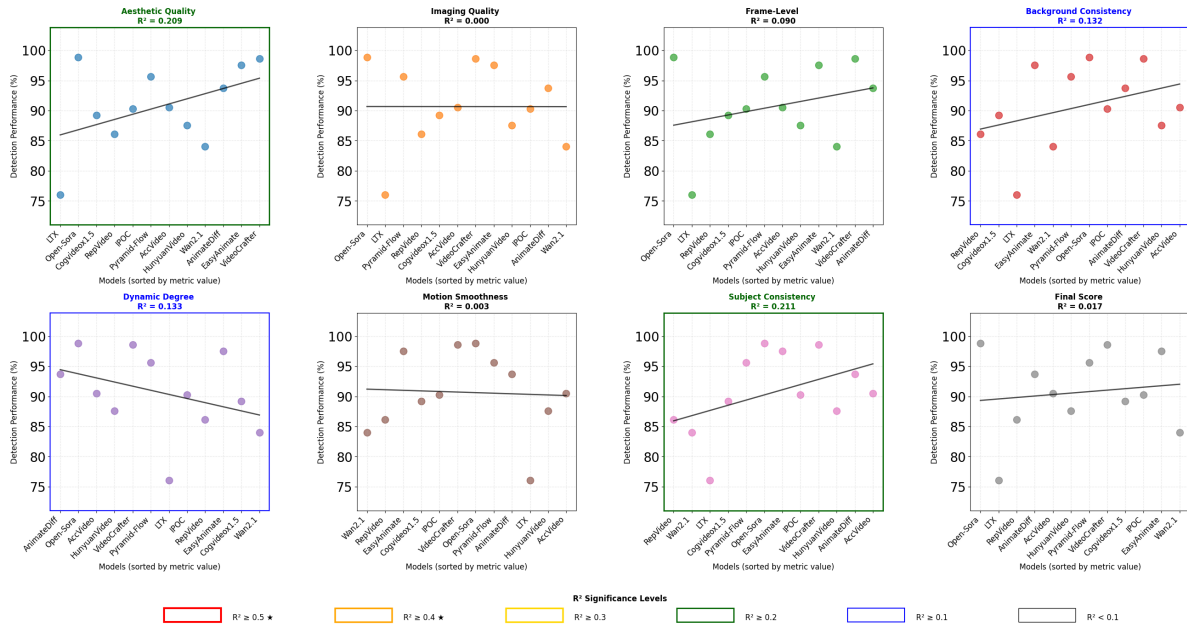


(a) Trained on Cogvidox1.5.

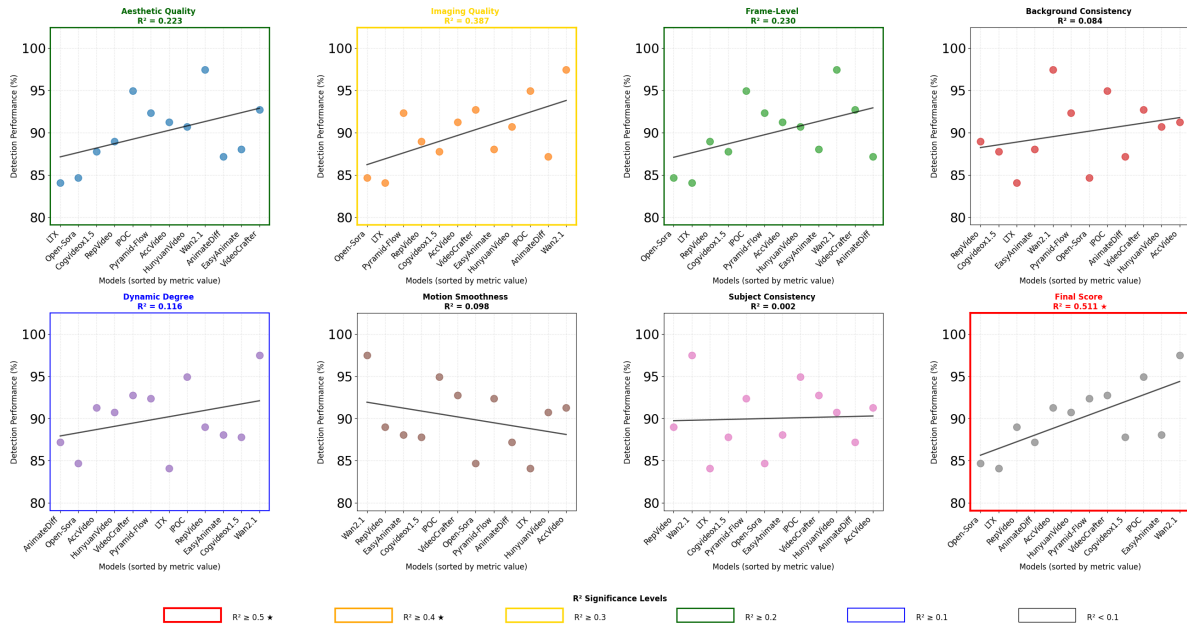


(b) Trained on IPOC.

Figure 66. Correlations between the performance of DeCoF and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

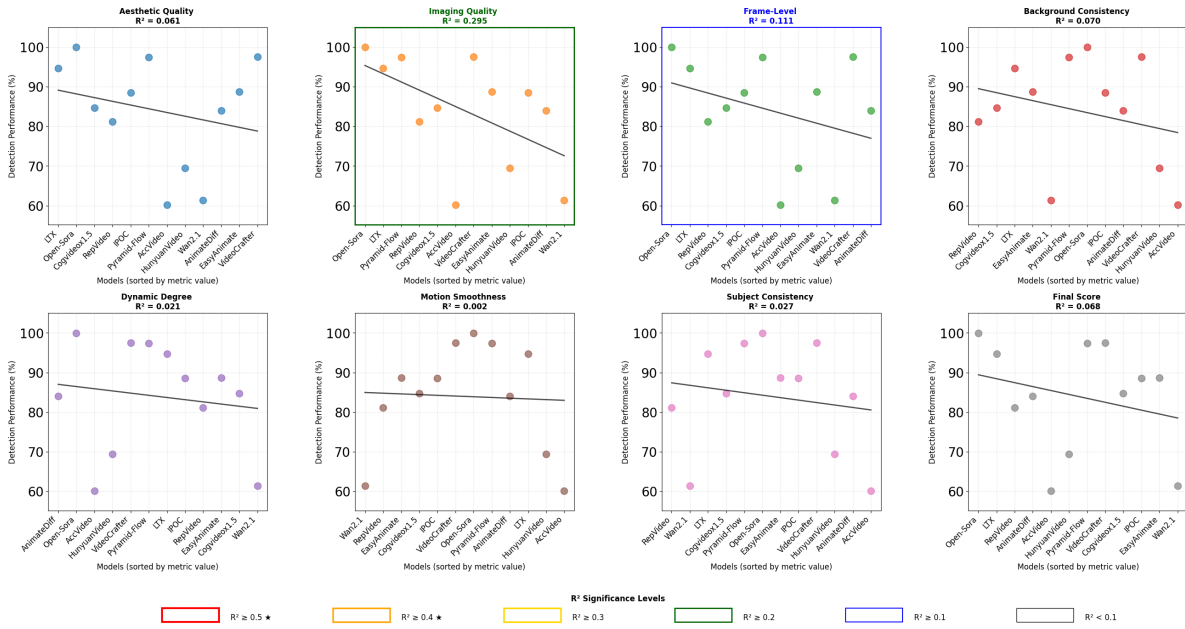


(a) Trained on EasyAnimate.

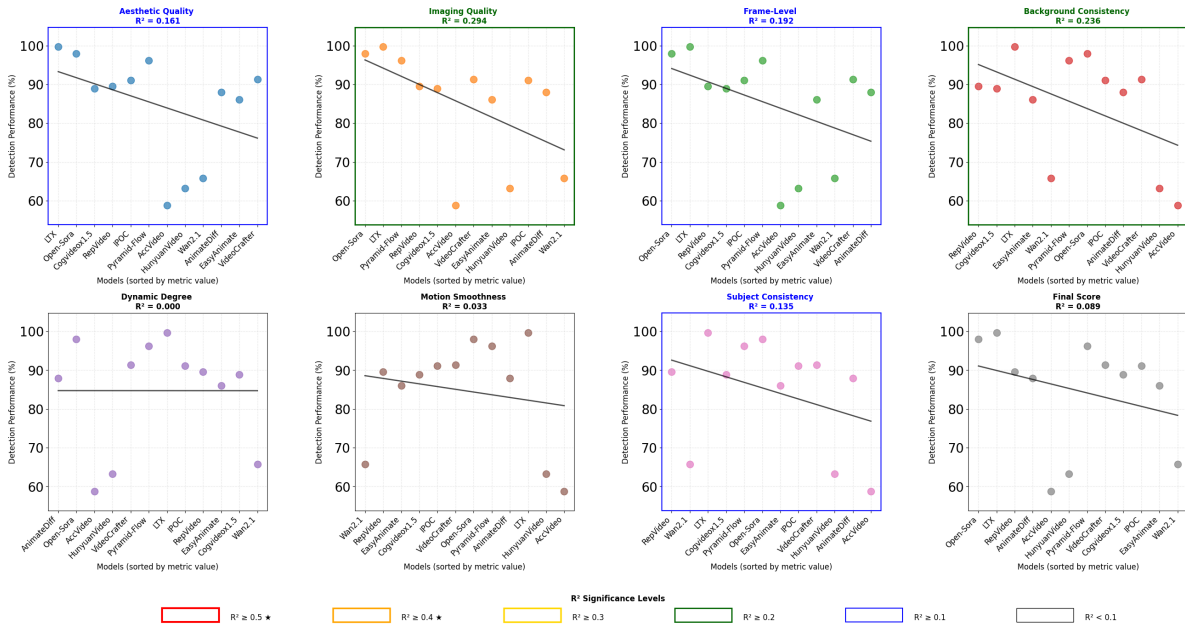


(b) Trained on Wan2.1.

Figure 67. Correlations between the performance of DeCoF and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

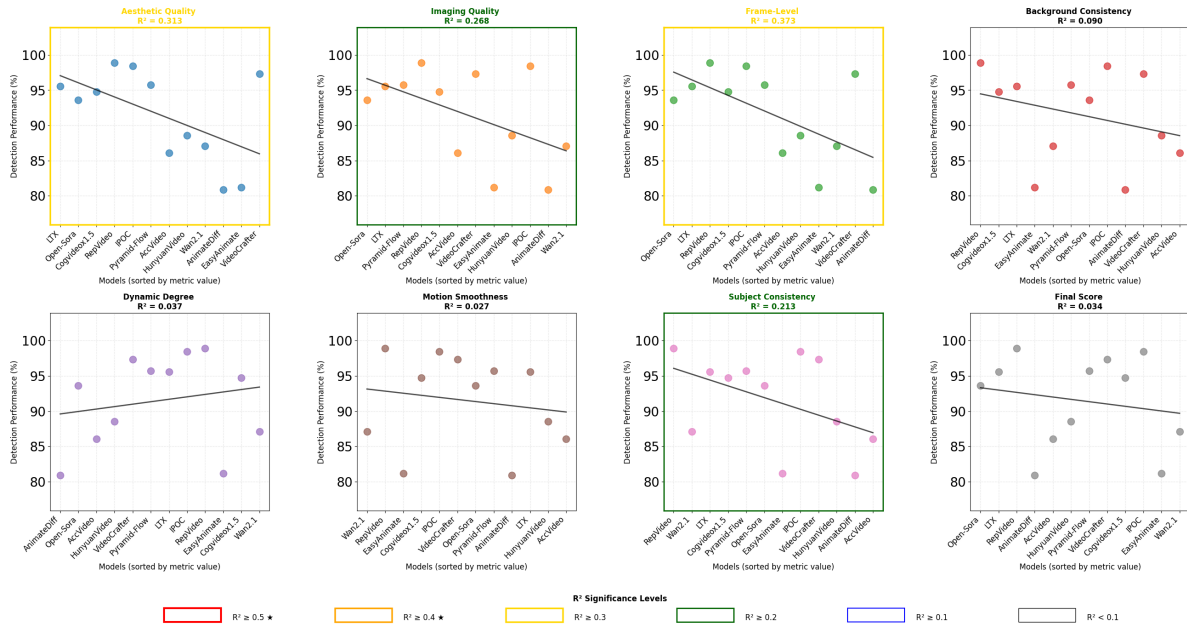


(a) Trained on Open-Sora.

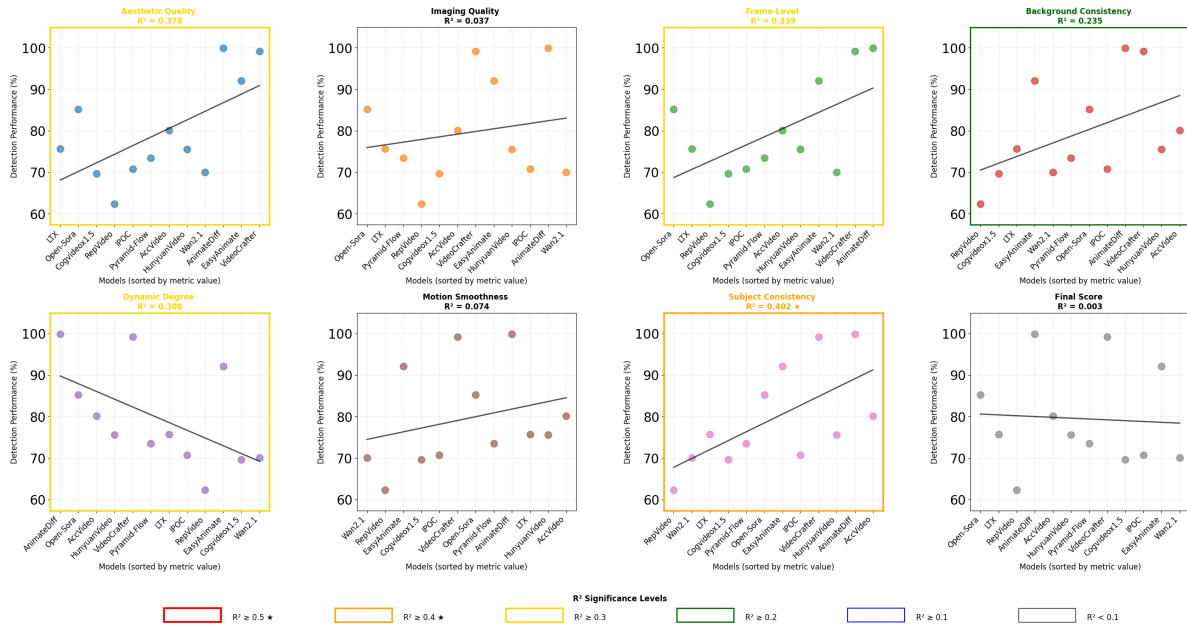


(b) Trained on LTX.

Figure 68. Correlations between the performance of UnivFD and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

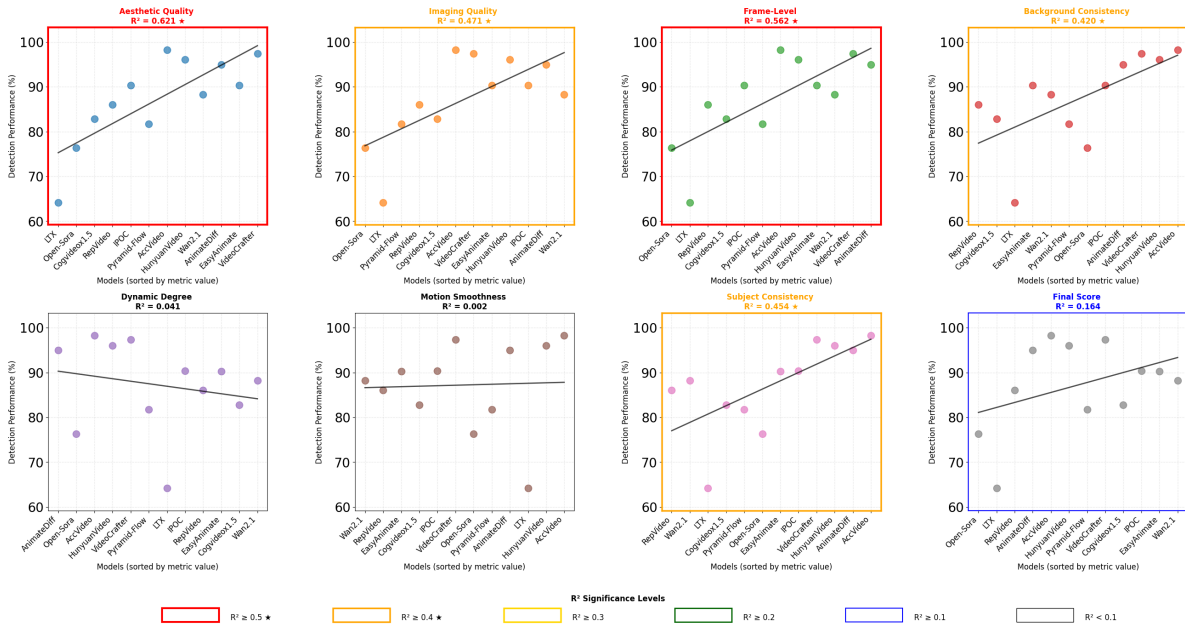


(a) Trained on RepVideo.

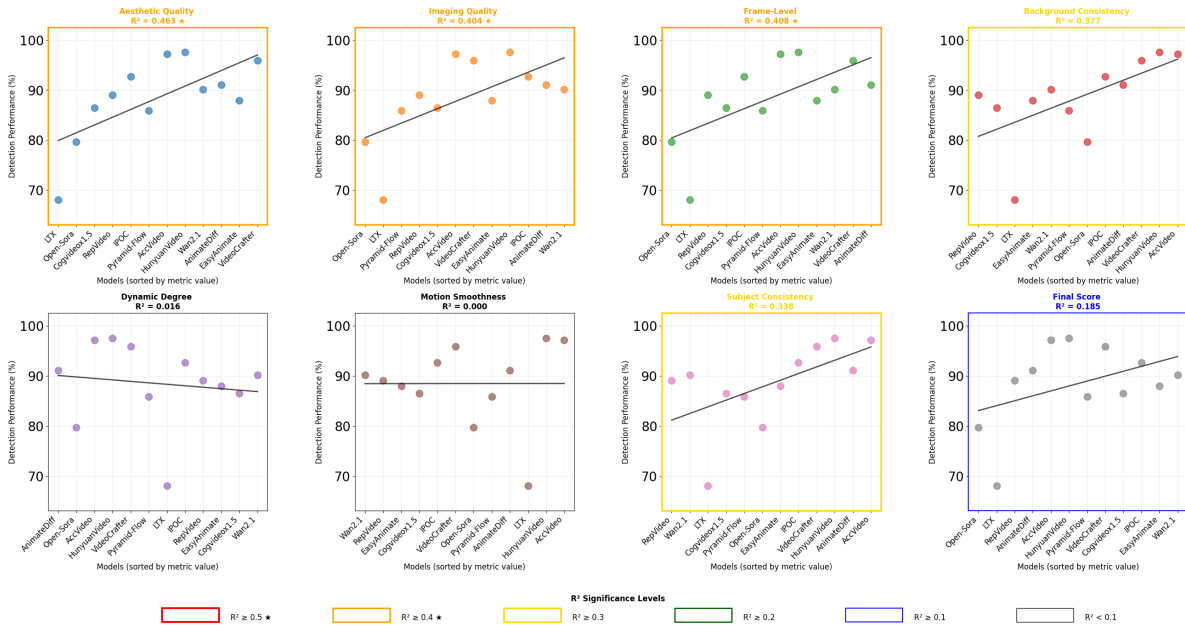


(b) Trained on AnimateDiff.

Figure 69. Correlations between the performance of UnivFD and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

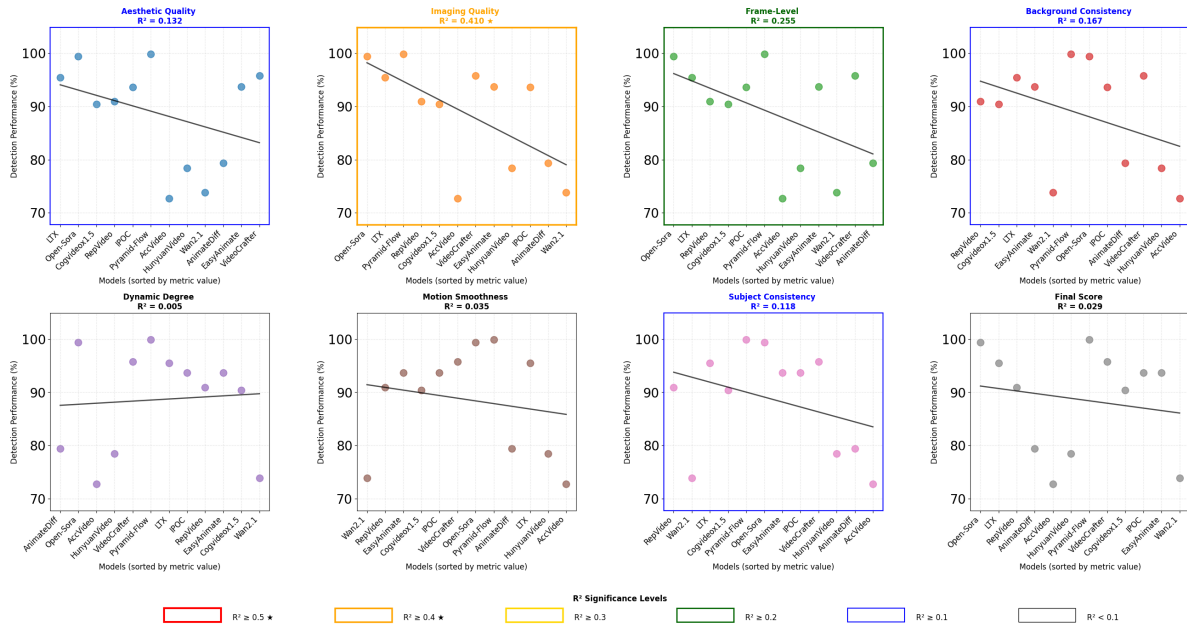


(a) Trained on AccVideo.

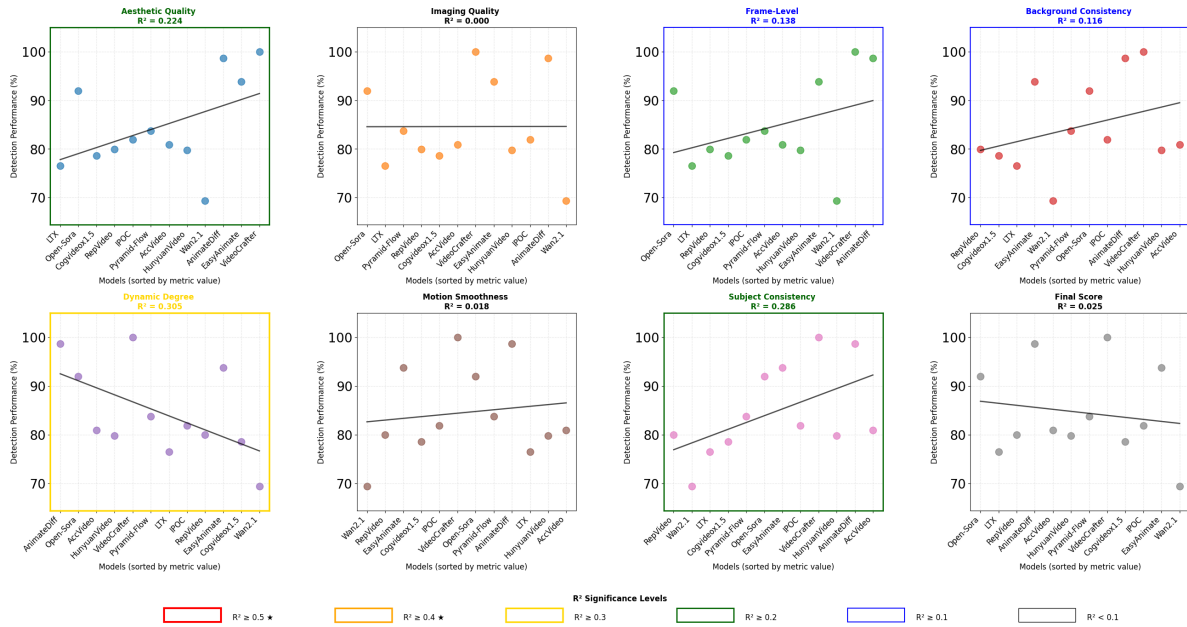


(b) Trained on HunyuanVideo.

Figure 70. Correlations between the performance of UnivFD and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

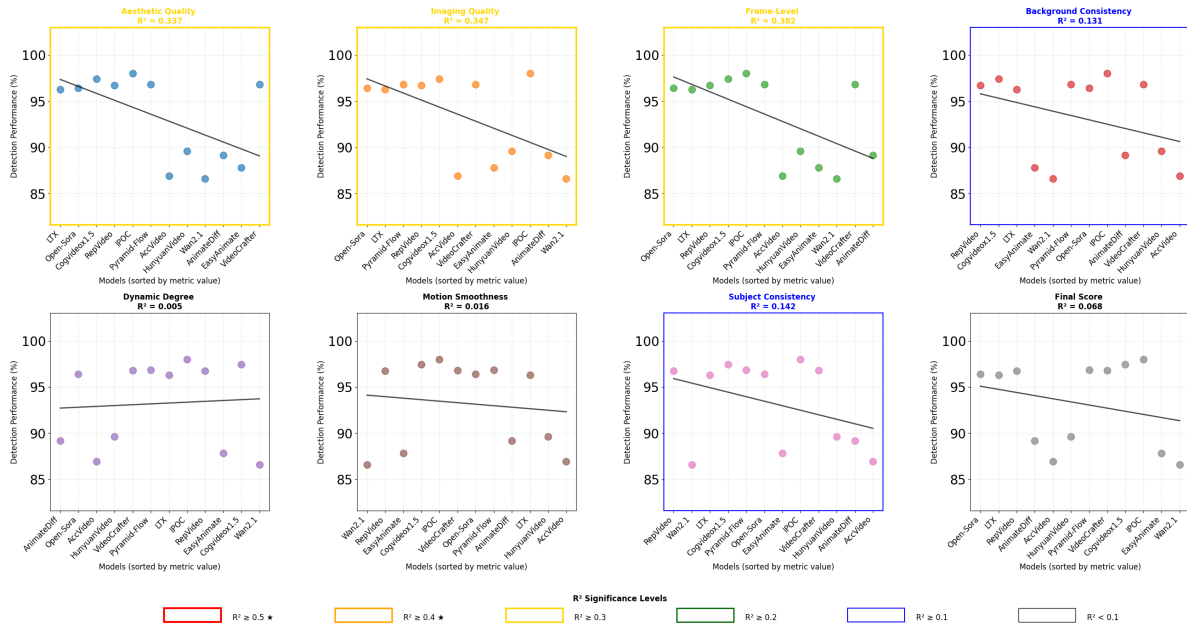


(a) Trained on Pyramid-Flow.

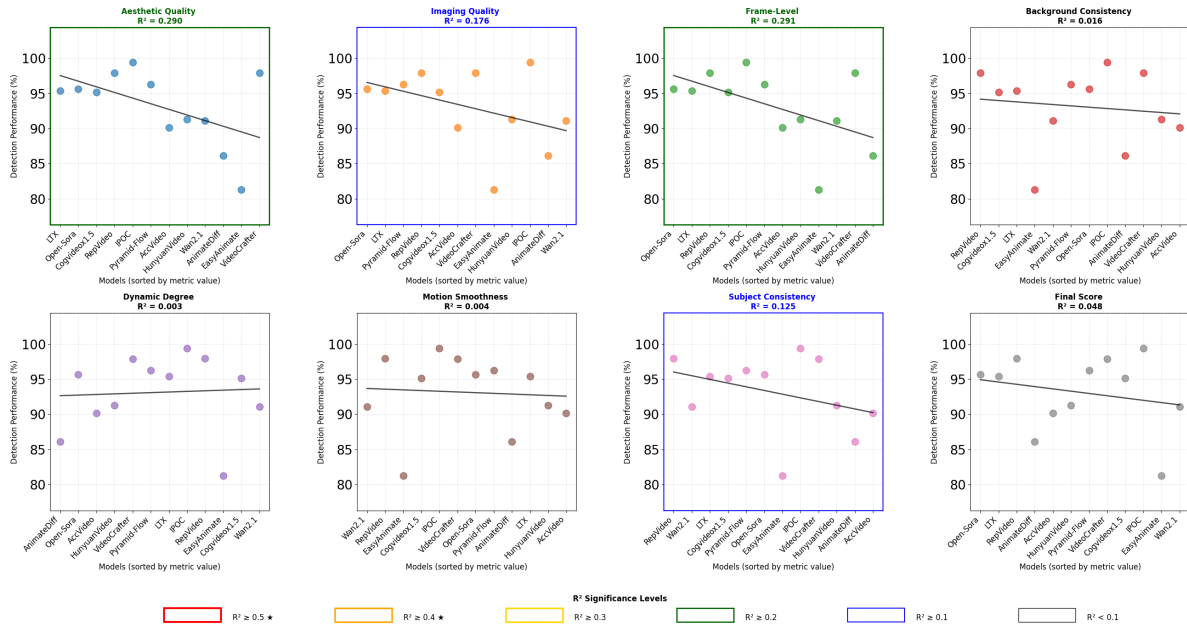


(b) Trained on VideoCrafter.

Figure 71. Correlations between the performance of UnivFD and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

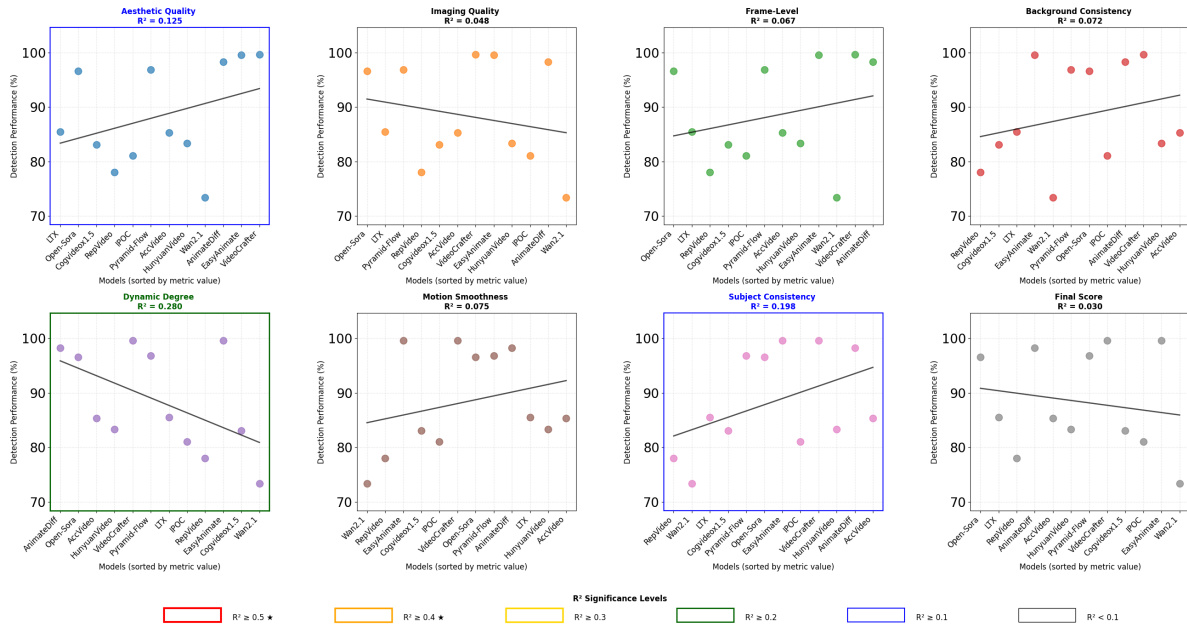


(a) Trained on Cogvidox1.5.

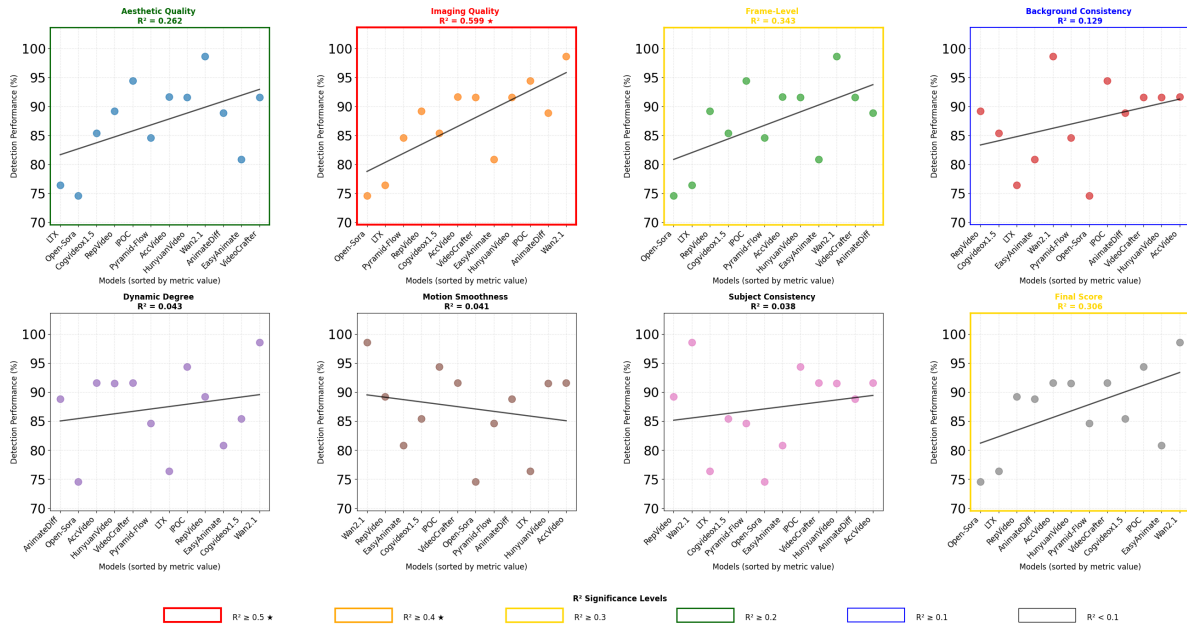


(b) Trained on IPOC.

Figure 72. Correlations between the performance of UnivFD and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

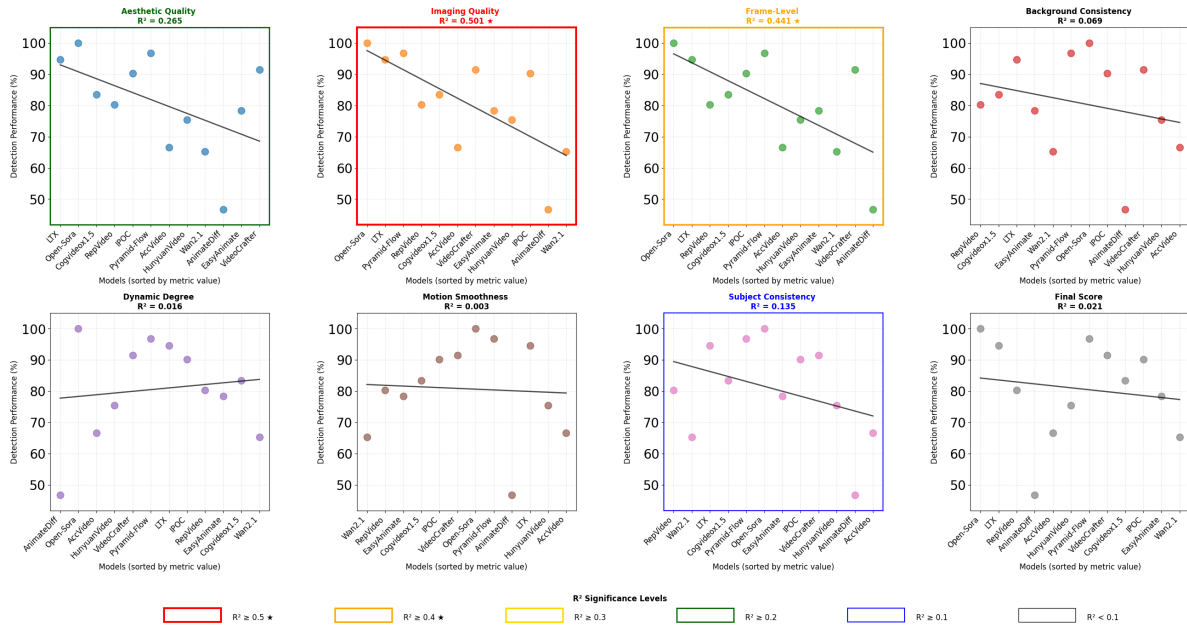


(a) Trained on EasyAnimate.

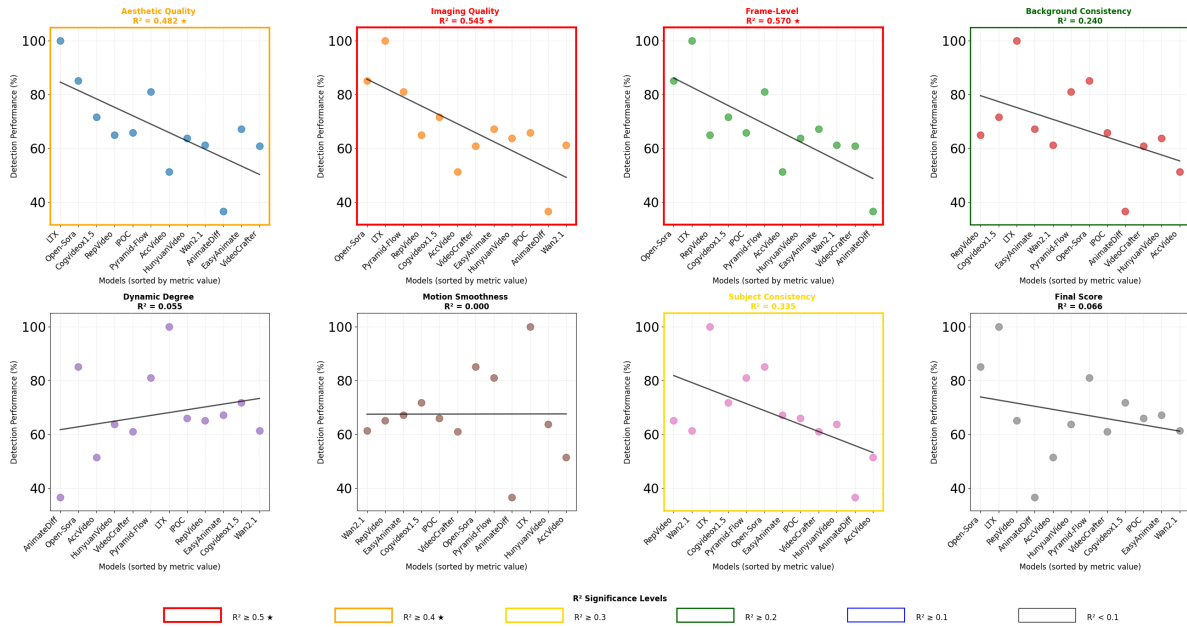


(b) Trained on Wan2.1.

Figure 73. Correlations between the performance of UnivFD and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

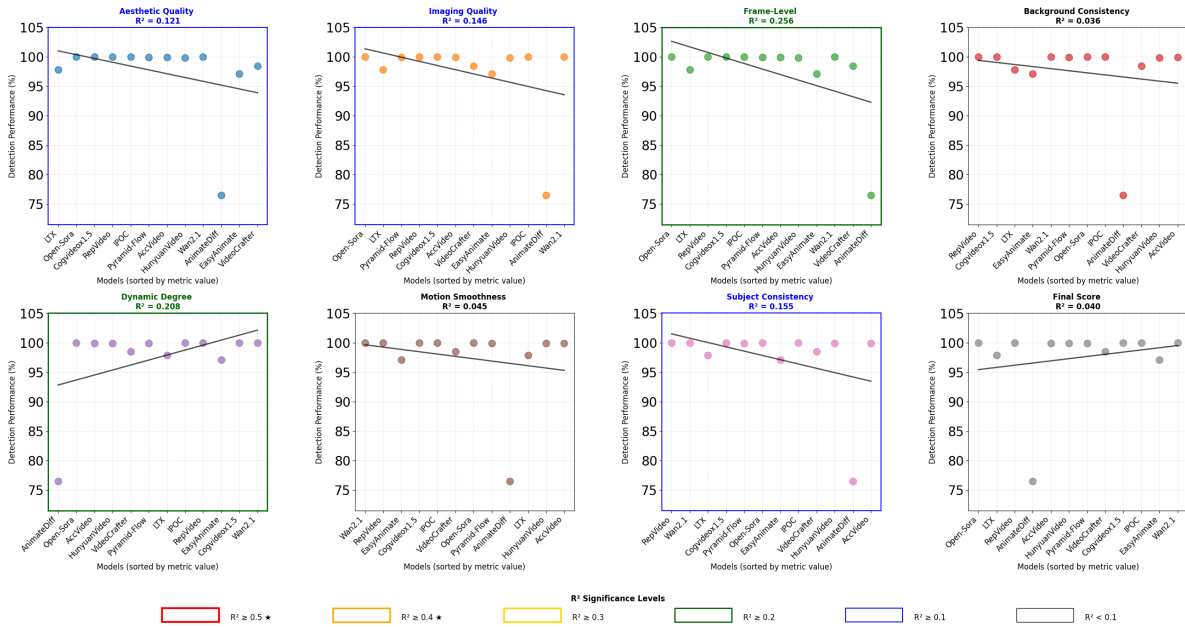


(a) Trained on Open-Sora.

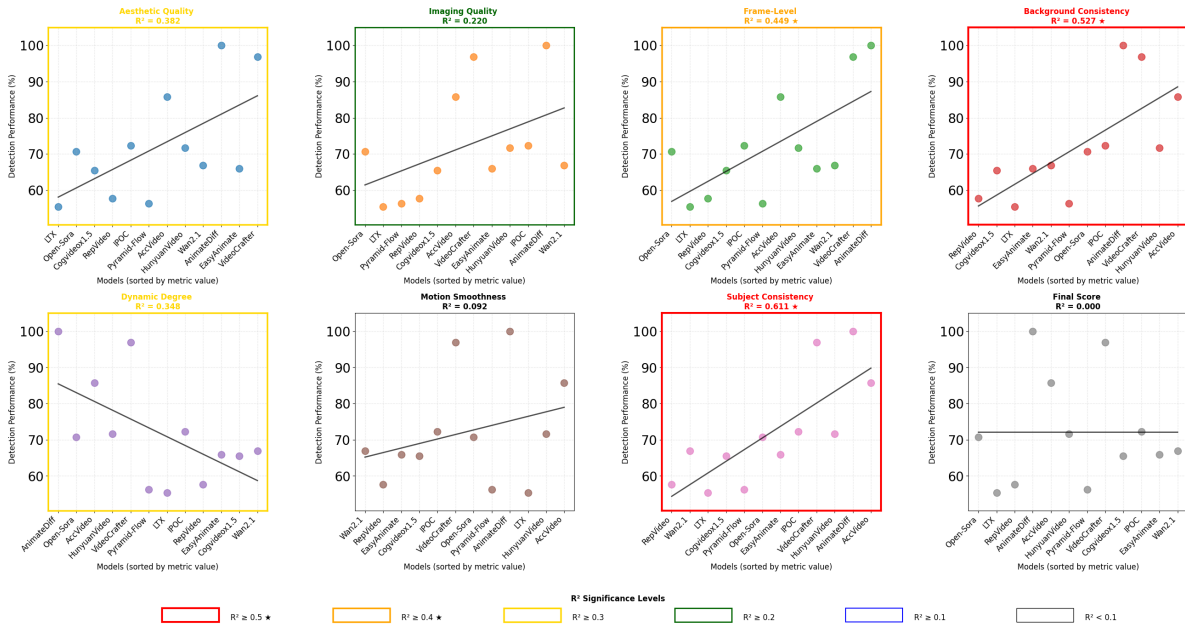


(b) Trained on LTX.

Figure 74. Correlations between the performance of I3D and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

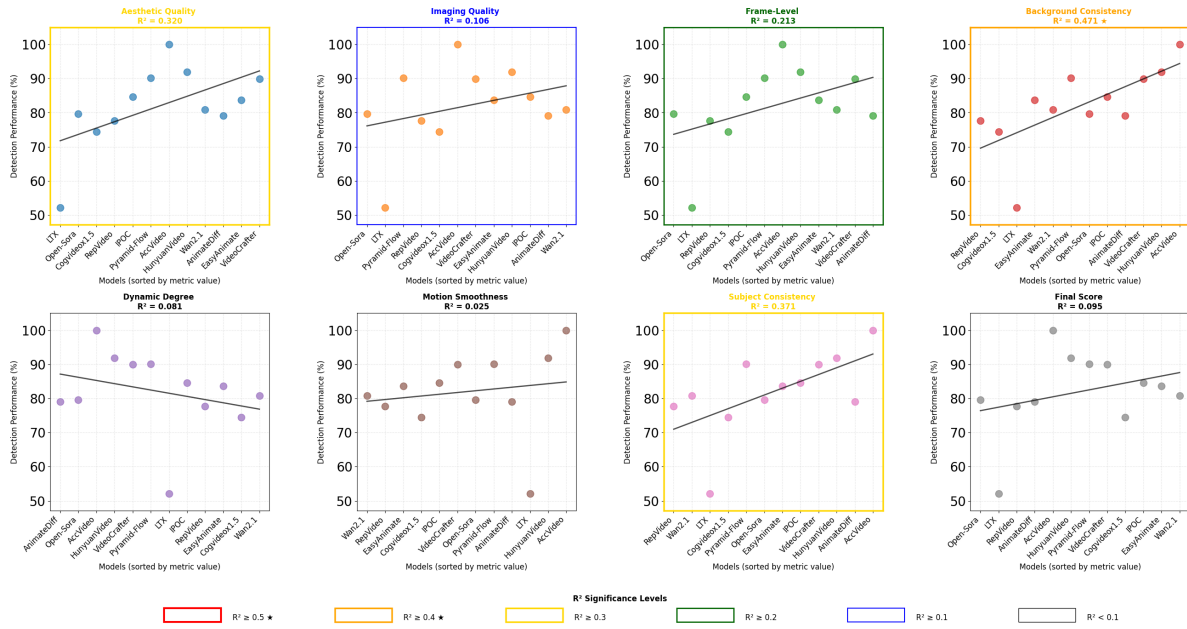


(a) Trained on RepVideo.

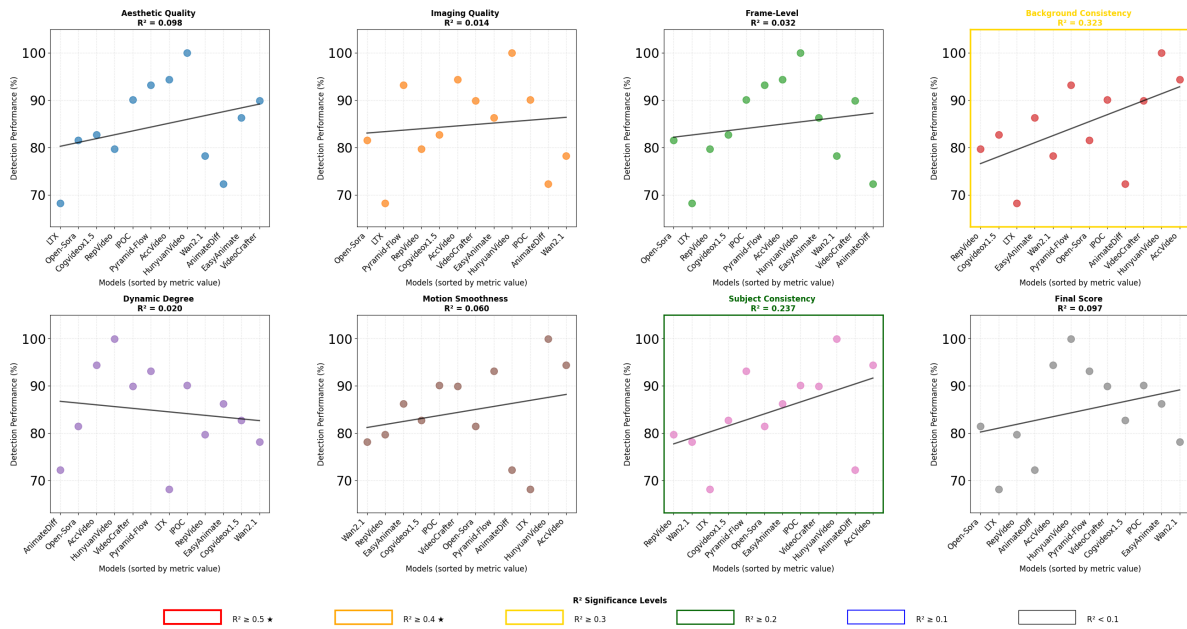


(b) Trained on AnimateDiff.

Figure 75. Correlations between the performance of I3D and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

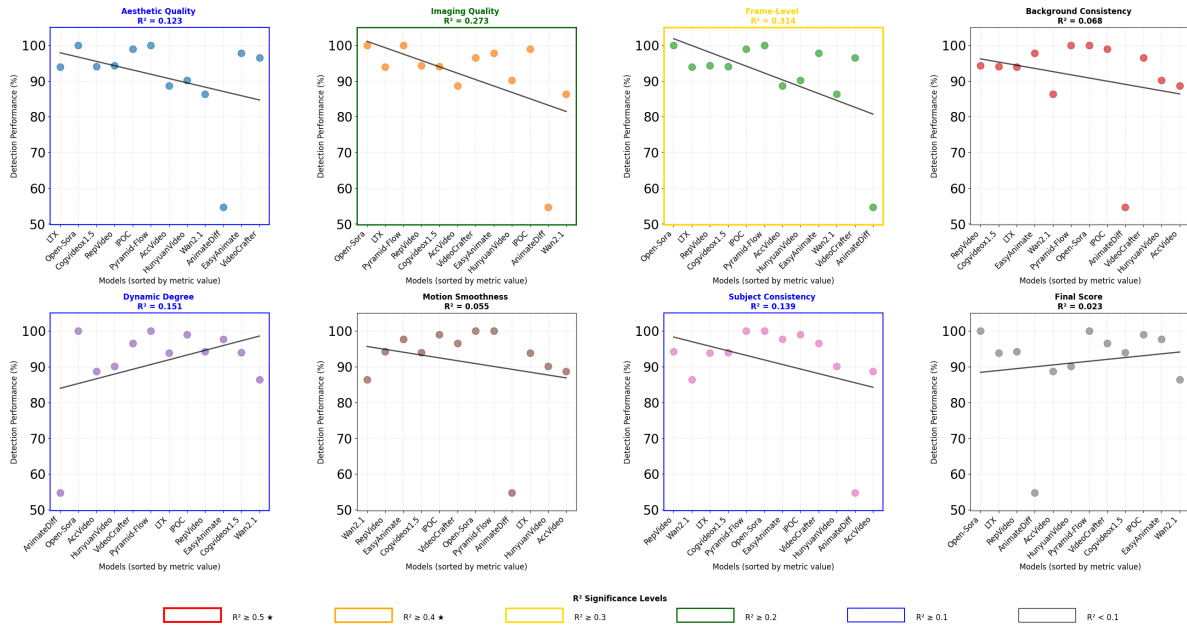


(a) Trained on AccVideo.

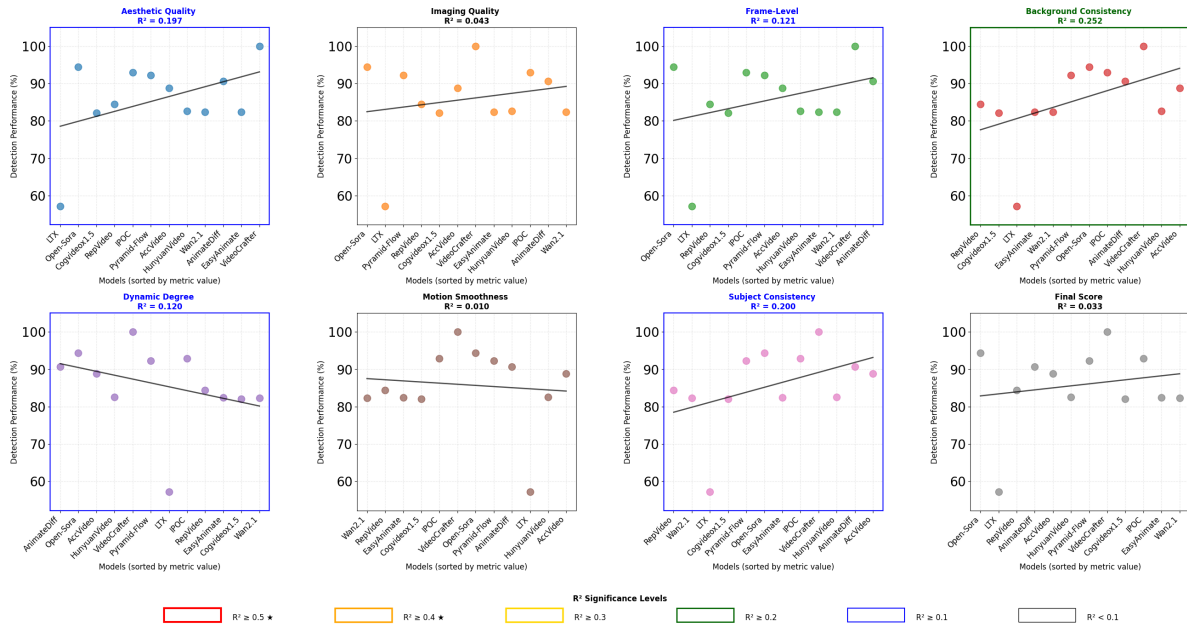


(b) Trained on HunyuanVideo.

Figure 76. Correlations between the performance of I3D and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

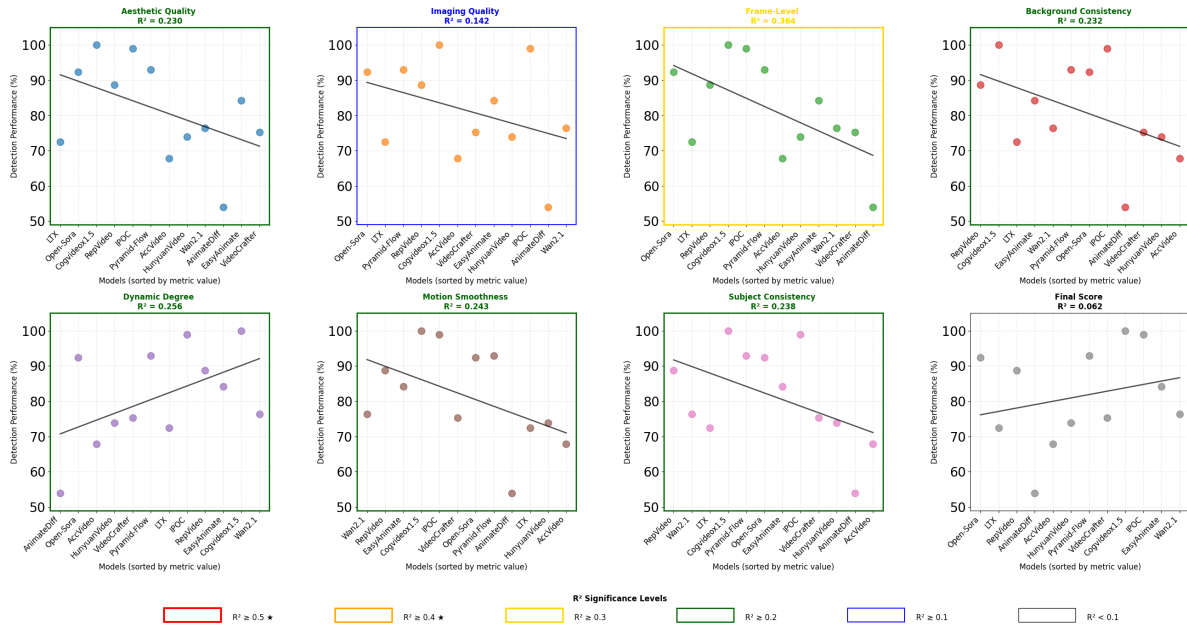


(a) Trained on Pyramid-Flow.

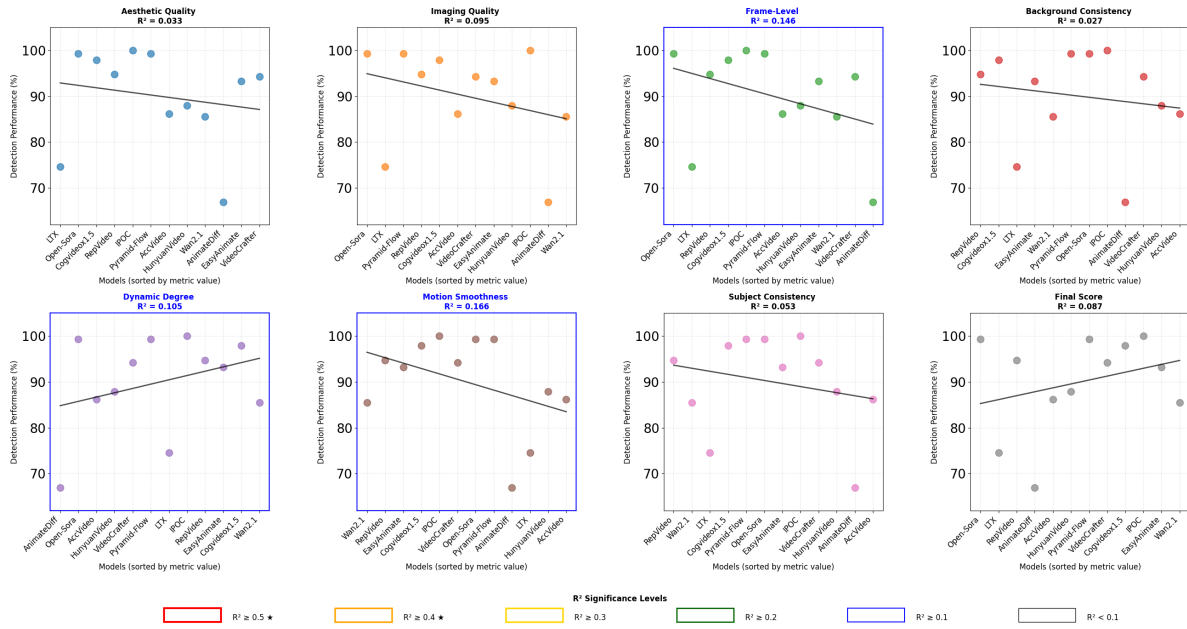


(b) Trained on VideoCrafter.

Figure 77. Correlations between the performance of I3D and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).

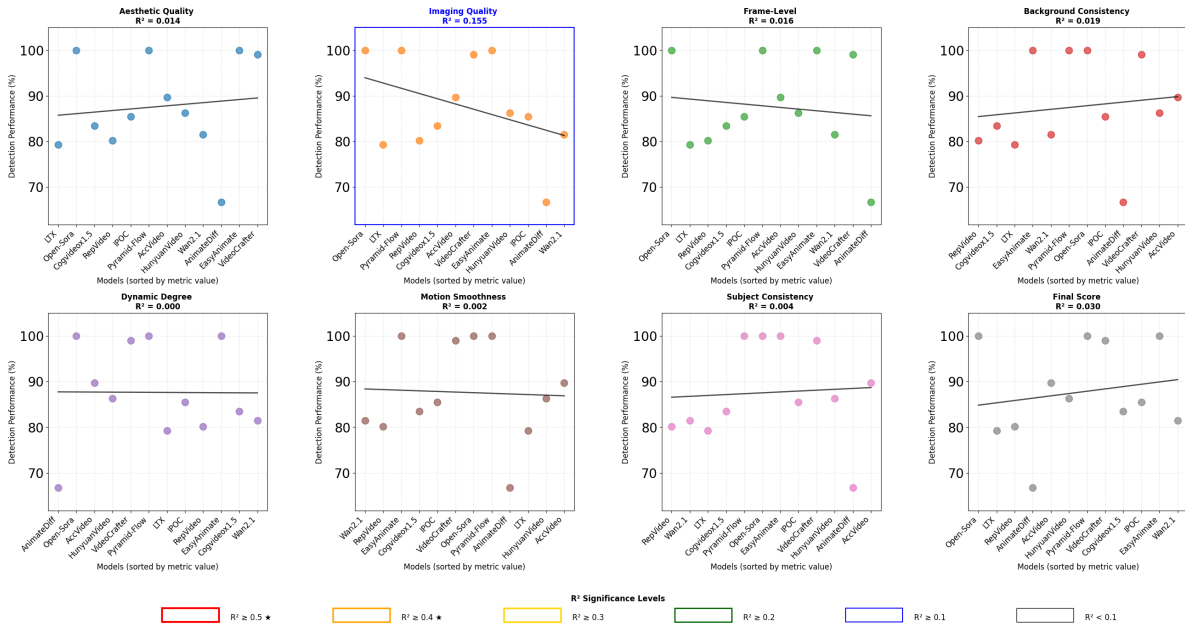


(a) Trained on Cogvidox1.5.

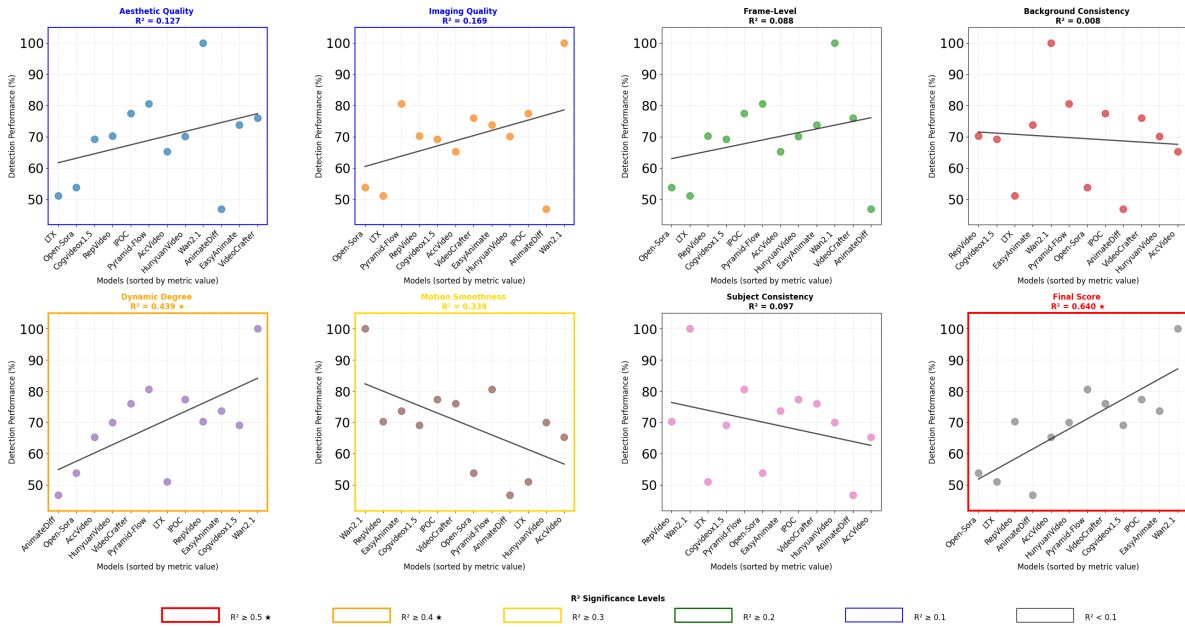


(b) Trained on IPOC.

Figure 78. Correlations between the performance of I3D and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).



(a) Trained on EasyAnimate.



(b) Trained on Wan2.1.

Figure 79. Correlations between the performance of I3D and quality metrics of video generation models. For each scatter plot, the metric value increases along the x-axis. The coefficient of determination (R^2) quantifies the strength of the linear relationship; a value closer to 1 indicates that the quality metric effectively explains the variance in detection performance (AUC).