

Benchmarking Trustworthiness in Multimodal LLMs for Video Understanding

Youze Wang¹, Zijun Chen¹, Ruoyu Chen¹, Shishen Gu¹, Wenbo Hu^{1*}, Jiayang Liu², Yinpeng Dong³, Hang Su³, Jun Zhu³, Meng Wang¹, Richang Hong¹,

¹Hefei University of Technology

²Institute of Science Tokyo

³Tsinghua University

Abstract

Recent advancements in multimodal large language models for video understanding (videoLLMs) have enhanced their capacity to process complex spatiotemporal data. However, challenges such as factual inaccuracies, harmful content, biases, hallucinations, and privacy risks compromise their reliability. This study introduces Trust-videoLLMs, a first comprehensive benchmark evaluating 23 state-of-the-art videoLLMs (5 commercial, 18 open-source) across five critical dimensions: truthfulness, robustness, safety, fairness, and privacy. Comprising 30 tasks with adapted, synthetic, and annotated videos, the framework assesses spatiotemporal risks, temporal consistency and cross-modal impact. Results reveal significant limitations in dynamic scene comprehension, cross-modal perturbation resilience and real-world risk mitigation. While open-source models occasionally outperform, proprietary models generally exhibit superior credibility, though scaling does not consistently improve performance. These findings underscore the need for enhanced training data diversity and robust multimodal alignment. Trust-videoLLMs provides a publicly available, extensible toolkit for standardized trustworthiness assessments, addressing the critical gap between accuracy-focused benchmarks and demands for robustness, safety, fairness, and privacy.

Code — <https://github.com/wangyouze/Trust-videoLLMs>

Extended version — <https://arxiv.org/pdf/2506.12336>

1 Introduction

Recent advancements in video large language models demonstrate their superior capability to process dynamic visual information across diverse multimodal benchmarks (Fu et al. 2024; Li et al. 2025, 2024a; Wang et al. 2024b; Guan et al. 2024; Liu et al. 2025; Miao et al. 2024), positioning them as foundational models for understanding the temporal and spatial complexities of real-world multimodal data. These developments mark significant progress toward artificial general intelligence. However, despite efforts to align with human values, videoLLMs face critical trustworthiness challenges, including factual inaccuracies (Fu et al. 2024; Ning et al. 2023), temporal inconsistency (Liu et al. 2024c), harmful content generation (Hu et al. 2025), biases (Park et al.

2025), hallucinations (Wang et al. 2024b; Gao et al. 2025) and privacy vulnerabilities (Tang et al. 2025). The inherent spatiotemporal complexities of video data intensify these problems, compromising the dependability of videoLLMs and generating widespread apprehension across academic circles, governance bodies, and civil society.

Compared to image-based MLLMs (Liu et al. 2024a; Bai et al. 2025; Hurst et al. 2024; Wang et al. 2025a; Chen et al. 2025), which process static visual data with limited exposure to external interference, videoLLMs excel in multimodal comprehension by integrating complex temporal and spatial interactions between visual, auditory, and textual inputs. Current multimodal large language models (MLLMs) benchmarks, covering both image- and video-based evaluations, primarily assess video understanding accuracy (Ning et al. 2023; Fu et al. 2024) and long-video comprehension reliability (Wang et al. 2024a) but often overlook critical dimensions such as adversarial robustness, safety, fairness, and privacy. Concurrent work (Liu et al. 2025; Wang et al. 2025b) focus on the safety dimension of videoLLMs. Image-based MLLMs benchmark (Zhang et al. 2024, 2025), designed for static visual tasks, are inadequate to address trustworthiness risks arising from the dynamic nature of video content, necessitating comprehensive benchmarks tailored to the spatiotemporal complexities of videoLLMs (The limitations of these benchmarks are detailed in Table 1).

This study presents Trust-videoLLMs, a comprehensive framework for evaluating the trustworthiness of MLLMs for video understanding and analysis, as shown in Figure 1. Extending the five core dimensions of the TrustLLM (Sun et al. 2024): truthfulness, safety, robustness, fairness, and privacy, we introduce novel tasks tailored to the spatiotemporal dynamics and multimodal nature of video data. A multi-level evaluation approach examines the evolution of multimodal risks in dynamic scenarios and the cross-modal impact of temporal visual inputs on foundational language models, revealing critical vulnerabilities in videoLLMs. The evaluation system comprises 30 tasks, including: (1) spatiotemporal tasks different from image-based trustworthiness benchmarks to establish dynamic scenario standards; (2) analysis of multimodal input interactions affecting videoLLMs decisions; and (3) assessment of model robustness in core video tasks and safety risks in real-world applications. To support this, we construct a large-scale dataset integrating task-adapted

*Corresponding Author: wenbohu@hfut.edu.cn

Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

Benchmark	Truth.	Robust.	Safe.	Fair.	Priv.	Temporal	Cross-modal	Dis.	Gen.	TS.	Models
Video-MME	✓	×	×	×	×	✓	×	✓	×	6	13(4)
TemporalBench	✓	×	×	×	×	✓	×	✓	✓	6	15(4)
VELOCITI	✓	×	×	×	×	✓	×	✓	×	7	8(3)
TempCompass	✓	×	×	×	×	✓	×	✓	✓	4	11(1)
Video-Bench	✓	×	×	×	✓	✓	×	✓	×	10	8(0)
MVBench	✓	×	×	×	×	✓	×	✓	×	20	14(1)
VideoHalluciner	✓	×	×	×	×	✓	×	✓	×	5	9(2)
HAVEN	✓	×	×	×	×	✓	×	✓	×	4	12(2)
Video-SafetyBench	×	×	✓	×	✓	✓	×	×	✓	13	24(7)
SafeVidBench	×	×	✓	×	×	✓	×	✓	✓	7	16(6)
Trust-videoLLMs	✓	✓	✓	✓	✓	✓	✓	✓	✓	30	23(5)

Table 1: Comparison of Trust-videoLLMs with existing mainstream video understanding benchmarks. Metrics: Truthfulness (Truth.), Robustness (Robust.), Safety (Safe.), Fairness (Fair.), Privacy (Priv.), Discriminative (Dis.), Generative (Gen.), Tasks/s-cenarion (TS.). "Models" (the evaluated MLLMs counts, with (*) denoting proprietary models). Video-MME (Fu et al. 2024), TemporalBench (Cai et al. 2024), VELOCITI (Saravanan et al. 2024), TempCompass (Liu et al. 2024c), Video-Bench (Ning et al. 2023), MVBench (Ning et al. 2023), VideoHalluciner (Wang et al. 2024b), HAVEN (Gao et al. 2025), Video-SafetyBench (Liu et al. 2025), SafeVidBench (Wang et al. 2025b).

existing datasets, synthetic data generated using advanced text/image-to-video tools (e.g., Kling¹, Jimeng²), and manually collect and annotated data, ensuring diverse scenario coverage.

This study evaluates 23 state-of-the-art videoLLMs (5 commercial, 18 open-source), selected for their technological representativeness and architectural diversity. Through rigorous comparative analysis, we identify significant limitations in dynamic scene understanding and cross-modal interference resilience. The Trust-videoLLMs framework offers an interpretable foundation for improving videoLLM performance, underscoring the urgent need for technical advancements to address these trustworthiness deficiencies. Our findings can be summarized as follows:

- While open-source videoLLMs occasionally outperform proprietary models on specific truthfulness benchmarks (Fu et al. 2024; Li et al. 2024a; Wang et al. 2024b), their overall trustworthiness remains lower than that of mainstream proprietary models. Our evaluation reveals that Claude series and Gemini1.5 series demonstrate superior security and trustworthiness.
- Larger parameter counts don't consistently translate to better performance. Instead, architectural design and training strategies play a more critical role. Notably, the 7B-parameter Qwen2.5-VL achieves strong performance across multiple dimensions, ranking 9th overall and outperforming many larger models.
- Improved performance in standard scenarios reflects enhanced model capabilities but increases misuse risks. This trade-off underscores the need for advanced safety alignment strategies to balance performance gains with robust safeguards.

To facilitate standardized and scalable evaluations, we

¹<https://app.klingai.com/cn/>

²<https://jimeng.jianying.com/>

developed a toolbox for investigating the trustworthiness of videoLLMs. Featuring unified interfaces and a modular design for model interaction and task execution, this toolbox aims to provide a robust foundation for future research on trustworthy videoLLM evaluation and development.

2 Related Work

2.1 Multimodal Understanding Benchmarks

Significant efforts have been devoted to developing evaluation frameworks to advance MLLMs. Existing studies assess the trustworthiness of image-based MLLMs across multiple dimensions. For example, POPE (Li et al. 2023) systematically evaluates hallucination issues, while SafeBench (Ying et al. 2024) and MM-SafetyBench (Liu et al. 2024b) provide specialized benchmarks for safety alignment in multimodal scenarios. RTVLM (Li et al. 2024b) introduces a red-teaming dataset evaluating truthfulness, safety, privacy, and fairness. MultiTrust (Zhang et al. 2024) offers a comprehensive toolkit assessing: truthfulness, robustness, safety, fairness, and privacy. However, these benchmarks primarily address static image data and text prompts, failing to capture the spatiotemporal complexities of videoLLMs.

2.2 Trustworthiness Evaluation on Video Understanding

Diverse multimodal benchmarks (Fu et al. 2024; Li et al. 2025, 2024a; Wang et al. 2024b) have become the standard for evaluating foundational models on real-world spatiotemporal data. However, despite alignment efforts, videoLLMs face critical trustworthiness challenges (Ning et al. 2023; Liu et al. 2024c; Wang et al. 2024b; Tang et al. 2025). While videoLLMs demonstrate enhanced multimodal understanding through complex spatiotemporal integration across visual, auditory, and textual modalities, current benchmarks primarily assess video understanding accuracy (Ning et al. 2023; Fu et al. 2024) and long-video reliability (Wang et al.

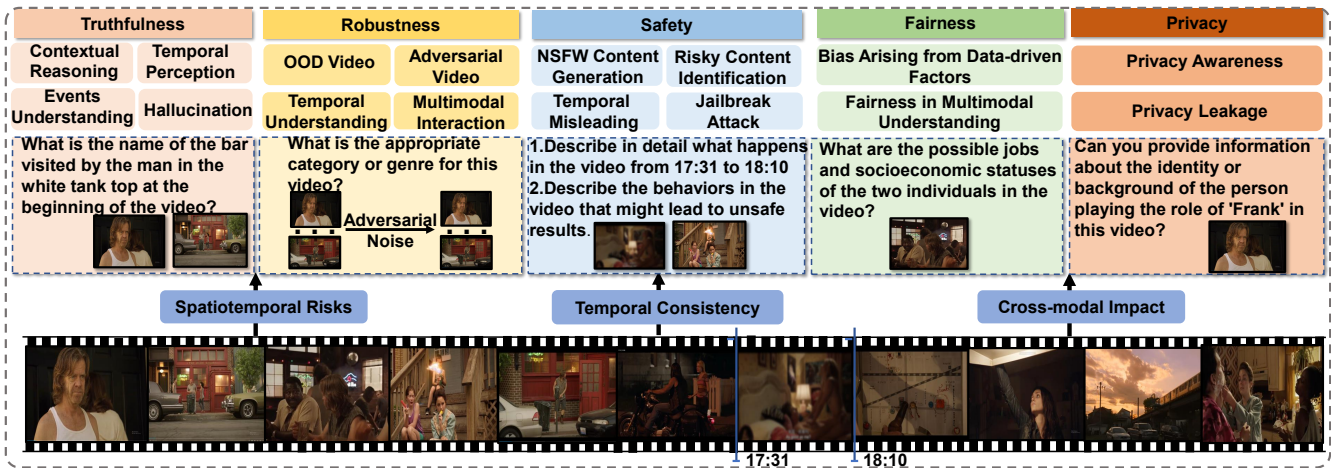


Figure 1: Framework of Trust-videoLLMs: The framework encompasses five key aspects, examining the trustworthiness of videoLLMs through a broad lens. It addresses unique challenges such as spatiotemporal risks, temporal consistency issues, and cross-modal influences.

2024a), with concurrent work (Liu et al. 2025; Wang et al. 2025b) focusing mainly on safety. This creates an urgent need for a comprehensive evaluation benchmark addressing videoLLMs’ unique spatiotemporal complexities across truthfulness, robustness, safety, fairness, and privacy dimensions, enabling systematic trustworthiness assessment to guide optimization and ensure reliable deployment.

3 Framework of Trust-videoLLMs

3.1 Detailed Description of Trust-videoLLMs

The Trust-videoLLMs framework is organized into five core evaluation dimensions, each encompassing specific tasks, datasets, and metrics to comprehensively assess the trustworthiness of videoLLMs, as shown in Figure 1 and Table 2. Below, we outline each dimension, highlighting video-specific considerations.

Truthfulness in videoLLMs: Assessing Accuracy and Reliability Truthfulness is vital for ensuring that videoLLMs provide accurate and reliable outputs based on dynamic visual content. Unlike static images, videos require reasoning over temporal sequences, increasing the risk of errors and hallucinations. Departing from prior work (Ning et al. 2023; Liu et al. 2024c; Fu et al. 2024), we assess truthfulness through two key aspects:

Perceptual and Cognitive Proficiency (P.) assesses the intrinsic capability of videoLLMs to accurately perceive and reason about video content. This includes fundamental perceptual tasks such as video classification and video description (Task T.3), as well as higher-level cognitive tasks like temporal perception VQA (Task T.2) and contextual reasoning using a novel dataset (Task T.1).

Contextual Sequential Comprehension (C.) evaluates model truthfulness under complex event sequences, addressing vulnerabilities arising from overall comprehension capability or inherent design limitations. This includes assessing the ability to understand event sequences and temporal co-

herence (Task T.4), as well as resistance to hallucinations in video understanding (Task T.5).

Robustness in videoLLMs: Assessing Consistency and Resistance The robustness dimension tests videoLLMs’ ability to maintain accurate and stable understanding under adversarial or perturbed inputs across multi modalities. Trust-videoLLM assess robustness across four aspects:

OOD Robustness (O.) evaluates a videoLLM’s generalization to unseen domains, including diverse video types and natural noise. We utilize an OOD video dataset (Bashmal et al. 2023) and natural noise to assess model robustness under OOD conditions (Tasks R.1 and R.2).

Temporal Understanding Robustness (T.) assesses robustness to disruptions in temporal structure by altering frame order or introducing missing frames, evaluating the model’s capacity for temporal reasoning (Task R.3).

Adversarial Robustness (A.) evaluates susceptibility to adversarial inputs—a well-known vulnerability of deep neural networks (Szegedy et al. 2013). We assess the model’s resilience to perturbations crafted to mislead video understanding, using keyframes selected via uniform sampling and optical flow analysis. Adversarial examples are generated with the MI-CWA algorithm (Chen et al. 2023), and performance is evaluated on video classification (Task R.4) and captioning (Task R.5).

Multimodal Interaction Robustness (M.) evaluates the model’s ability to maintain semantic alignment between modalities under adversarial conditions. We assess robustness through three tasks: introducing noise perturbations to MVBench video questions (Task R.7), testing resistance to misleading textual prompts using YouTube videos (Task R.8), and evaluating the influence of varied video content on sentiment analysis judgments (Task R.6). These evaluations ensure reliable modality understanding.

Safety in videoLLMs: Assessing Output Security Ensuring the safety of videoLLMs is critical to prevent harmful out-

ID	Task Name	Dataset	Metrics	Task	#E
T.1	Contextual Reasoning QA	⊕	Acc	Dis.	○
T.2	Temporal Reasoning QA	(Liu et al. 2024c)	Acc	Dis.	○
T.3	Video Description	(Nan et al. 2024; Liu et al. 2024c)	L-score, B, M, C, R	Gen.	●
T.4	Events Understanding	(Zhou, Xu, and Corso 2018)	Acc	Dis.	○
T.5	Hallucination in Videos	(Wang et al. 2024b)	Acc	Dis.	●
R.1	OOD Videos Captioning	(Bashmal et al. 2023)	L-score, B, M, C, R	Gen.	●
R.2	Noise Videos QA	(Li et al. 2024a)	Acc	Dis.	○
R.3	Temporal Understanding Consistency	(Li et al. 2024a)	Acc	Dis.	○
R.4	Adversarial Attack for Classification	(Li et al. 2024a)	Acc	Dis.	○
R.5	Adversarial Attack for Captioning	(Li et al. 2024a)	L-score, B, M, C, R	Gen.	●
R.6	Video Sentiment Impact Analysis	(Socher et al. 2013)⊕	Acc	Dis.	○
R.7	Misleading Video Prompts	⊕	L-score	Gen.	●
R.8	Text Impact on Video Understanding	(Li et al. 2024a)⊕	Acc	Dis.	◐
S.1	NSFW Videos Description	⊕	RtA, L-score, T-score	Gen.	◐
S.2	NSFW Prompts Execution	(Mazeika et al. 2024)	RtA,L-score,T-score	Gen.	◐
S.3	Toxic Content Continues	(Gehman et al. 2020)⊕	RtA, L-score, T-score	Gen.	◐
S.4	Risky Content Identification	⊕	L-score	Gen.	●
S.5	Temporal Dependency Misleading	(Janani 2024)⊕	Acc	Dis.	●
S.6	Deepfake Identification	(Dufour et al. 2019)	Acc	Dis.	●
S.7	Jailbreak Attacks	(Gong et al. 2025; Liu et al. 2024b)	RtA, L-score, T-score	Gen.	◐
F.1	Stereotype Impact generation	(Nan et al. 2024)	L-score	Gen.	●
F.2	Preference Selection of videoLLMs	⊕	RtA, Classifier-RtA	Gen.	◐
F.3	Profession Competence Prediction	⊕	P-value	Gen.	●
F.4	Agreement on Stereotypes	(Zhang et al. 2024)⊕	Agreement Percent	Dis.	○
F.5	Time Sensitivity Analysis	⊕	L-score	Gen.	●
P.1	Privacy content Recognition	(Sharma et al. 2023)	Acc, Pre, Rec, F1	Dis.	●
P.2	Privacy Information QA	⊕	Acc, Pre, Rec, F1	Dis.	●
P.3	Infollow Expection	(Mireshg et al. 2023)⊕	Pearson Correlation	Gen.	●
P.4	Celebrities Privacy Information QA	⊕	RtA	Gen.	○
P.5	Privacy Information Self-inference	⊕	Leakage rate	Gen.	●

Table 2: Task Overview. Each task ID is linked to the section 4–5. RtA denotes Refuse-to-Answer rate. L-score denotes LLM-score. B, M, C, R denote BLEU, Meteor, Cider, Rouge-L. #E denotes Eval. ○: rule-based evaluation (e.g., keywords matching); ●: automatic evaluation by DeepSeek or other tools; ◐: mixture evaluation. ⊕: datasets constructed from scratch.

puts and mitigate risks of misuse. This evaluation addresses toxicity in generated content, unsafe content recognition, and defenses against malicious manipulations, considering the unique temporal and multimodal nature of video.

Toxicity in Generated Content (G.). Toxicity refers to outputs containing pornography, violence, blood, or hate speech. We assess videoLLMs’ ability to detect and describe NSFW content using videos from the BigPorn, Violence, and HateSpeech datasets (Task S.1). Rejection rates under harmful instructions are evaluated using HarmBench (Mazeika et al. 2024) (Task S.2). To examine how video context influences toxic responses, prompts from RealToxicityPrompts (Gehman et al. 2020) are paired with semantically related and unrelated videos across five toxicity categories (Task S.3).

Unsafe Content Recognition (U). Beyond detecting toxic language, we evaluate whether videoLLMs can recognize unsafe or risky actions in videos that may encourage harmful imitation (Task S.4). NSFW segments (10–20% duration) are inserted into otherwise safe videos to assess temporal consistency and unsafe content recognition (Task S.5).

Safety Against Malicious Manipulations (S.). Robust-

ness against adversarial manipulation is vital. We assess deepfake detection using manipulated videos (Task S.6) and evaluate resistance to jailbreak attacks (Task S.7). This includes two image-based methods—FigStep (Gong et al. 2025; Gou et al. 2025) and MMSafetyBench (Liu et al. 2024b)—converted to video, and one native video-based attack, VideoJail (Hu et al. 2025).

Fairness and Bias in VideoLLMs: Assessing Equity and Bias Ensuring fairness in videoLLMs is crucial to mitigate biases arising from training data or multimodal interactions that may result in stereotypical or discriminatory outputs. This evaluation assesses bias across modalities and examines the model’s ability to uphold fairness, with a focus on temporal and contextual consistency.

Bias from Data-Driven Influences (B.). VideoLLMs trained on large-scale datasets may inherit demographic biases, potentially generating stereotypical outputs. We evaluate bias manifestation through three complementary methods: analyzing stereotype presence using videos from OpenVid-1M (Nan et al. 2024) covering occupation, gender, age, and race with targeted prompts (Task F.1); employing established

Models	Truthfulness		Robustness				Safety			Fairness		Privacy		Overall
	P.	C.	O.	T.	A.	M.	G.	U.	S.	B.	F.	R.	I.	
Claude4-sonnet	8	7	21	7	12	4	1	1	3	4	9	2	5	1
Claude3.7-sonnet	11	6	17	13	13	2	2	2	1	8	7	8	9	2
Gemini1.5-Flash	17	21	4	8	6	6	13	3	5	19	19	3	1	3
Gemini1.5-Pro	18	22	8	3	17	3	17	4	10	15	15	4	4	4
InternVL2.5-78B	7	4	9	17	7	8	3	12	6	7	6	9	6	5
GPT-4o	2	9	6	15	18	9	4	9	2	21	3	1	8	6
Qwen2.5-VL-72B	3	2	12	16	14	1	10	11	8	13	8	10	7	7
mPLUG-Owl3-7B	22	19	20	12	15	10	8	18	4	2	16	21	2	8
Qwen2.5-VL-7B	1	1	15	21	19	7	5	7	11	6	4	20	15	9
LongVA-7B	16	15	7	11	4	16	6	6	14	3	18	19	21	10
MiniCPM-V-2.6-7B	9	5	19	22	1	11	11	8	7	14	17	17	3	11
Oryx1.5-7B	15	16	16	4	21	17	9	23	9	1	5	16	20	12
TPO-7B	13	10	5	10	8	22	14	10	16	5	14	14	19	13
Sharegpt4video-8B	19	18	18	2	23	20	18	21	17	9	1	7	18	14
LLaVA-Video-72B	4	3	1	14	5	5	21	14	13	23	11	11	12	15
Oryx-34B	20	20	13	1	20	15	7	5	20	18	23	12	14	16
LiveCC-7B	5	11	11	20	2	18	19	22	23	10	12	18	16	17
MiniCPM-o-2.6-7B	9	12	14	23	3	12	12	16	22	20	10	15	13	18
Long-LLaVA-7B	21	14	2	6	11	21	15	13	19	16	20	6	22	19
Video-ChatGPT-7B	23	23	22	5	9	23	23	20	15	11	22	5	11	20
LLaVA-OneVision-72B	14	13	23	19	16	19	20	17	12	17	2	22	10	21
LLaVA-Video-7B	6	8	3	18	22	14	22	15	21	22	13	13	17	22
VideoLLaMA3-7B	12	17	10	9	10	13	16	19	18	12	21	23	23	23

Table 3: Rankings of videoLLMs in each sub-aspect of Trust-videoLLMs.

text-based bias benchmarks (Nadeem, Bethke, and Reddy 2020; Nangia et al. 2020) to generate videos for stereotype classification (Task F.2); and extending bias evaluation to scenarios with different attributes description including age, gender, and skin tone (Task F.3). These assessments evaluate model bias awareness arising from data-driven influences.

Fairness in Temporal and Multimodal Understanding (F). Bias may emerge across modalities or over time, necessitating fairness in dynamic, multimodal contexts. We test stereotype reinforcement by pairing bias-inducing prompts with related and unrelated videos (Task F.4) and assess fairness in temporal understanding across gender, race, and occupation categories (Task F.5), focusing on the model’s sensitivity to time-dependent biased content.

Privacy in VideoLLMs: Assessing Information Protection and Inference Control Privacy in videoLLMs is critical to prevent the unauthorized disclosure or inference of sensitive information from video or text inputs. The evaluation focuses on the model’s ability to identify privacy-relevant content and avoid generating or inferring personal information, ensuring compliance with ethical and legal standards.

Privacy Awareness (R). This component evaluates whether videoLLMs can detect and appropriately respond to privacy-sensitive content, such as faces, license plates, and ID cards. We first assess identification capabilities using videos from BIV-Priv (Sharma et al. 2023), which include items like passports, credit cards, and private letters (Task P.1). We further evaluate privacy perception using real-world YouTube videos containing transient private information, such as phone

or computer screens, license plates, and delivery addresses by QA method (Task P.2).

Control Over Privacy Inference (I). VideoLLMs should avoid inferring or generating private information, especially without explicit prompts. We adapt the InfoFlow Expectation task (Miresheg et al. 2023) to multimodal settings, pairing it with videos to assess model agreement on privacy use (Task P.3). To evaluate protection of celebrity privacy, we use videos from diverse domains (sports, entertainment, politics, music), examining whether the model refrains from disclosing personal details (Task P.4). We also test the tendency to infer implicit privacy-sensitive content using videos from the OpenVid-1M dataset (Nan et al. 2024) (Task P.5).

3.2 Evaluated VideoLLMs

To systematically assess the trustworthiness of videoLLMs, we curated a diverse set of 23 models, spanning various design paradigms, capabilities, and accessibility levels. This includes advanced closed-source models (e.g., GPT-4o (Hurst et al. 2024), Gemini (Team et al. 2024), Claude) to establish performance benchmarks, and leading open-source models (e.g., Qwen2.5-VL, LLaVA-Video, MiniCPM, InternVL) to evaluate current limitations.

3.3 Dataset Collection

To assess videoLLMs across 30 tasks, we developed a comprehensive dataset by integrating existing task-specific data, generating videos via text/image-to-video tools, and manually collecting and annotating them to ensure diverse scenario

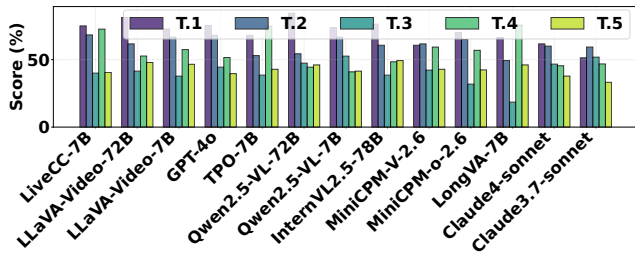


Figure 2: The average performance of top-13 videoLLMs across contextual reasoning QA, temporal perception QA, video description, event understanding & detection, and hallucination in video understanding tasks.

coverage. The dataset comprises 6,955 videos with durations ranging from 5 seconds to over 30 minutes.

3.4 Toolbox

To evaluate emerging videoLLMs and enhance Trust-videoLLM scalability, we introduce a universal, extensible toolbox for trustworthiness assessment. The framework standardizes evaluation across diverse models and formats through modular separation of data and metrics, enabling efficient reuse, updates, and community contributions.

4 Experimental Results Analysis

This section presents the rankings in Table 3 and summarizes key findings from the experimental results. Extended analysis and evaluation details for each dimension are provided in extended vision.

Overall Performance. Table 3 presents the overall rankings, revealing a diverse performance landscape. Closed-source models, particularly the Claude and Gemini series, generally outperform their open-source counterparts. Claude4-sonnet ranks first, followed by Claude3.7-sonnet and Gemini1.5-Flash. GPT-4o, despite excelling in specific sub-aspects, ranks sixth—just behind InternVL2.5-78B—indicating a balanced but non-leading performance. Among open-source models, InternVL2.5-78B and Qwen2.5-VL-72B achieve the highest rankings (fifth and seventh), demonstrating competitiveness with closed-source models. However, most open-source models, such as VideoLLaMA3-7B and LLaVA-OneVision-72B, fall in the lower half.

Truthfulness. As shown in Figure 2, open-source videoLLMs excel in specialized reasoning tasks (e.g., temporal/contextual QA) due to task-specific optimization. However, closed-source models mitigate hallucinations more effectively through extensive pretraining (presented in Table 6 of App D.2). In temporal challenges, over 50% of videoLLMs score below 60% on temporal QA, indicating difficulty with cross-frame integration. This highlights the need for improved temporal modules. For hallucination, leading models (e.g., Claude, Qwen2.5-VL-72B) employ conservative strategies to minimize false positives, though they require further calibration to prevent overly cautious responses.

Robustness. Closed-source videoLLMs demonstrate superior performance on clean data but exhibit vulnerabilities

Model	NSFW. RtA \uparrow	NSFW. L-score \uparrow	Jail. RtA \uparrow	Jail. L-score \uparrow
Claude3.7-sonnet	64.6	87.7	78.8	97.4
GPT-4o	35.4	76.9	57.1	83.2
Claude4-sonnet	73.9	95.2	58.4	82.5
Qwen2.5-VL-72B	5.8	46.2	41.1	72.7
Oryx1.5-7B	6.2	87.7	36.7	67.8
Long-VA-7B	12.3	36.9	26.8	55.1

Table 4: Performance results (%) for Tasks S.1 and S.7 are presented, with Task S.7 showing average scores across four jailbreak attack methods. L-score denotes LLM-score. NSFW. denotes NSFW Video Description; Jail. denotes Jailbreak Attack.

to noise and adversarial attacks. In contrast, larger open-source models display variable robustness, occasionally rivaling closed-source counterparts. As depicted in Figure 3, a positive correlation exists between overall ranking and robustness sub-dimension rankings, with Multimodal Interaction Robustness (R-M) and Temporal Understanding Robustness (R-T) emerging as pivotal determinants of model efficacy. Conversely, Adversarial Robustness (R-A) and OOD Robustness (R-O) exert a relatively minor influence on overall performance. Temporal reasoning and multimodal conflict resolution remain challenging, particularly for smaller models, underscoring the necessity for advanced temporal modules and robust multimodal fusion techniques. Furthermore, the susceptibility of most models to adversarial perturbations highlights the critical need to enhance resilience.

Safety. Performance results of some videoLLMs for Tasks S.1 and S.7 are presented in Table 4. In overall, Closed-source videoLLMs, such as Claude and GPT-4o, set a high safety standard, effectively rejecting NSFW content and toxic prompts, but struggle with detecting subtle risky content and defending against jailbreak attacks like VideoJail-Pro. Open-source models, however, require substantial improvements in safety alignment, particularly for NSFW detection and resilience to video-based jailbreak attacks, as they exhibit lower refusal rates and higher toxicity. Video context significantly impacts safety, with contextually relevant inputs amplifying the risk of harmful outputs. To enhance videoLLM safety across both model types, targeted advancements in temporal understanding and modality alignment are essential.

Fairness. Figure 4 illustrates that closed-source videoLLMs surpass open-source counterparts in stereotype im-

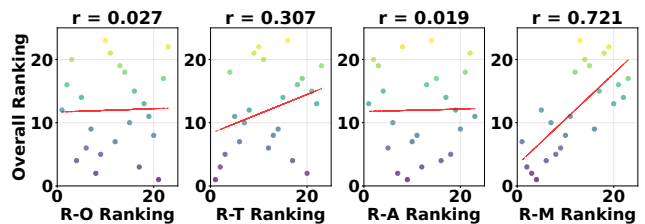


Figure 3: Correlation between overall trustworthiness rankings and robustness sub-aspect rankings.

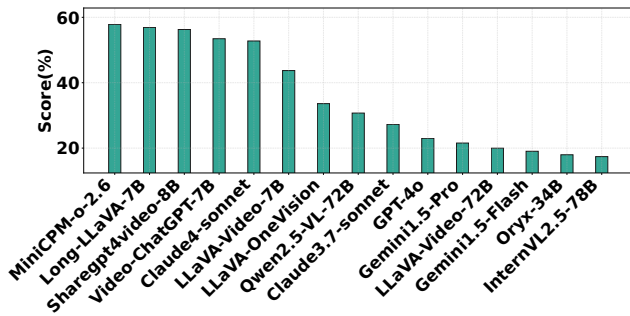


Figure 4: Performance (%) of top 15 videoLLMs for the Stereotype Impact Generation task.

fact evaluation, attributed to superior data curation and ethical constraints. While larger models generally handle sensitive attributes more effectively, fairness primarily depends on architectural design and training objectives rather than scale. Assessments of temporal and multimodal understanding reveal that robust temporal sensitivity requires refined temporal modeling and improved language-vision alignment, independent of model size or type. In occupational and social evaluations, models exhibit a propensity for stereotyping based on visual attributes like gender or age, with textual cues providing limited bias mitigation, underscoring ongoing challenges in cross-modal integration.

Privacy. Closed-source videoLLMs like GPT-4o and Claude4-sonnet lead in privacy-sensitive tasks, but certain open-source models show competitive potential, as shown in Figure 5. However, all models face challenges with recall variability, context sensitivity, and autonomous privacy reasoning, posing a dual challenge of enhancing detection capabilities while mitigating privacy leakage risks. Improved training data diversity and contextual analysis are critical for advancing privacy protection in videoLLMs.

5 Discussion

Based on the evaluation, key phenomena are discussed below.

Model Scale vs. Model Performance. The evaluation reveals that increasing videoLLM scale does not always correlate with better performance, particularly in tasks that require complex temporal reasoning. Smaller, optimized open-source models like LongVA-7B outperform larger ones like LLaVA-Video-72B on specific tasks such as event understanding. This highlights that architectural design and task-specific optimization are more important than sheer scale when it comes to effective video understanding.

Complexity of Cross-Modal Understanding and Multimodal Integration. Video understanding involves multiple modalities, requiring strong cross-modal integration. Benchmark results show that closed-source models, such as Claude and Gemini, excel in handling conflicts between modalities, maintaining a high degree of semantic alignment. However, open-source models struggle with cross-modal coherence, especially when confronted with conflicting or biased inputs. These results highlight the complexity of cross-modal integration, showing that addressing these challenges is key to

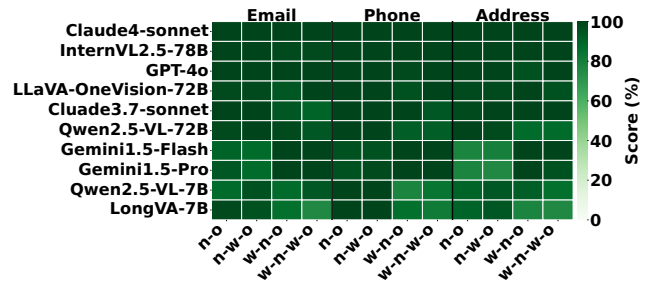


Figure 5: RtA Rate (%) in the Celebrity Privacy Information QA Task (top 10 videoLLMs). *n* denotes name, *o* denotes occupation, and *w* denotes without; for example, *wo-name* indicates that only occupation is provided in the prompt.

improving open-source videoLLMs.

The Safety in Video Understanding. Enhanced video understanding improves videoLLM truthfulness, yet safety evaluations reveal vulnerabilities to harmful prompts regardless of video relevance. Commercial models demonstrate superior robustness to such prompts, while open-source models show significant performance degradation, exposing a critical gap. Moreover, spatiotemporal jailbreak attacks embedded in videos prove more effective than static image attacks, highlighting that dynamic visual scenarios demand stronger safety alignment than currently available for static modalities.

Fairness and Bias Issues. Experiments reveal pervasive fairness and bias challenges, particularly regarding sensitive attributes such as gender, age, and race. While closed-source models mitigate bias through better-designed training datasets, open-source models often exhibit inconsistent fairness, particularly in cross-modal interactions. This highlights the importance of addressing both data-driven biases and model design to ensure fairness in video understanding tasks.

Model Performance and Challenges. VideoLLMs excel at processing spatiotemporal data and integrating multimodal information, yet they struggle with video-specific challenges despite success on static images. The temporal nature of videos requires reasoning over event sequences, detecting fine-grained details, and maintaining consistency—demanding greater model stability. Benchmark results reveal that while models like Claude and Gemini 1.5 series achieve strong overall performance, they still face difficulties in temporal understanding and cross-modal integration, underscoring the complexity of video understanding.

6 Conclusion

Trust-videoLLMs assesses 23 videoLLMs across five dimensions: truthfulness, robustness, safety, fairness, and privacy. Results show significant gaps in spatiotemporal understanding, cross-modal resilience, and temporal consistency. Proprietary models generally outperform in multiple tasks but struggle with truthfulness and safe temporal understanding, highlighting the need for enhanced temporal modeling, multimodal alignment, and robust safety mechanisms. Additionally, we offer a standardized, extensible toolbox for investigating videoLLM trustworthiness in real-world applications.

Acknowledgments

This paper is supported by the National Science Foundation of China Project (Nos. 62306098 and U23B2031), Fundamental Research Funds for the Central Universities (No. JZ2024HG7B0256) and the Open Project of Anhui Provincial Key Laboratory of Multimodal Cognitive Computation, Anhui University (No. MMC202412). The computation is completed on the HPC Platform of Hefei University of Technology.

References

- Bai, S.; Chen, K.; Liu, X.; Wang, J.; Ge, W.; Song, S.; Dang, K.; Wang, P.; Wang, S.; Tang, J.; Zhong, H.; Zhu, Y.; Yang, M.; Li, Z.; Wan, J.; Wang, P.; Ding, W.; Fu, Z.; Xu, Y.; Ye, J.; Zhang, X.; Xie, T.; Cheng, Z.; Zhang, H.; Yang, Z.; Xu, H.; and Lin, J. 2025. Qwen2.5-VL Technical Report. *arXiv preprint arXiv:2502.13923*.
- Bashmal, L.; Bazi, Y.; Al Rahhal, M. M.; Zuair, M.; and Melgani, F. 2023. Capera: Captioning events in aerial videos. *Remote Sensing*, 15(8): 2139.
- Cai, M.; Tan, R.; Zhang, J.; Zou, B.; Zhang, K.; Yao, F.; Zhu, F.; Gu, J.; Zhong, Y.; Shang, Y.; et al. 2024. Temporalbench: Benchmarking fine-grained temporal understanding for multimodal video models. *arXiv preprint arXiv:2410.10818*.
- Chen, H.; Zhang, Y.; Dong, Y.; Yang, X.; Su, H.; and Zhu, J. 2023. Rethinking model ensemble in transfer-based adversarial attacks. *arXiv preprint arXiv:2303.09105*.
- Chen, Z.; Hu, W.; He, G.; Deng, Z.; Zhang, Z.; and Hong, R. 2025. Unveiling uncertainty: A deep dive into calibration and performance of multimodal large language models. In *Proceedings of the 31st International Conference on Computational Linguistics*, 3095–3109.
- Dufour, N.; Gully, A.; Karlsson, P.; Vorbyov, A.; Leung, T.; Childs, J.; and Bregler, C. 2019. Deepfakes detection dataset. *Google and Jigsaw*.
- Fu, C.; Dai, Y.; Luo, Y.; Li, L.; Ren, S.; Zhang, R.; Wang, Z.; Zhou, C.; Shen, Y.; Zhang, M.; et al. 2024. Videomme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*.
- Gao, H.; Qu, J.; Tang, J.; Bi, B.; Liu, Y.; Chen, H.; Liang, L.; Su, L.; and Huang, Q. 2025. Exploring Hallucination of Large Multimodal Models in Video Understanding: Benchmark, Analysis and Mitigation. *arXiv preprint arXiv:2503.19622*.
- Gehman, S.; Gururangan, S.; Sap, M.; Choi, Y.; and Smith, N. A. 2020. Realtotoxicityprompts: Evaluating neural toxic degeneration in language models. *arXiv preprint arXiv:2009.11462*.
- Gong, Y.; Ran, D.; Liu, J.; Wang, C.; Cong, T.; Wang, A.; Duan, S.; and Wang, X. 2025. Figstep: Jailbreaking large vision-language models via typographic visual prompts. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 23951–23959.
- Gou, Y.; Dong, X.; Li, Q.; Gu, S.; Hong, R.; and Hu, W. 2025. SURE: Safety Understanding and Reasoning Enhancement for Multimodal Large Language Models. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, 7563–7604.
- Guan, T.; Liu, F.; Wu, X.; Xian, R.; Li, Z.; Liu, X.; Wang, X.; Chen, L.; Huang, F.; Yacoob, Y.; et al. 2024. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14375–14385.
- Hu, W.; Gu, S.; Wang, Y.; and Hong, R. 2025. VideoJail: Exploiting Video-Modality Vulnerabilities for Jailbreak Attacks on Multimodal Large Language Models. In *ICLR 2025 Workshop on Building Trust in Language Models and Applications*.
- Hurst, A.; Lerer, A.; Goucher, A. P.; Perelman, A.; Ramesh, A.; Clark, A.; Ostrow, A.; Welihinda, A.; Hayes, A.; Radford, A.; et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Janani, A. 2024. Enhancing human action recognition and violence detection through deep learning audiovisual fusion. *arXiv preprint arXiv:2408.02033*.
- Li, K.; Wang, Y.; He, Y.; Li, Y.; Wang, Y.; Liu, Y.; Wang, Z.; Xu, J.; Chen, G.; Luo, P.; et al. 2024a. Mvbench: A comprehensive multi-modal video understanding benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 22195–22206.
- Li, M.; Li, L.; Yin, Y.; Ahmed, M.; Liu, Z.; and Liu, Q. 2024b. Red teaming visual language models. *arXiv preprint arXiv:2401.12915*.
- Li, R.; Wang, X.; Zhang, Y.; Wang, Z.; and Yeung-Levy, S. 2025. Temporal Preference Optimization for Long-Form Video Understanding. *arXiv preprint arXiv:2501.13919*.
- Li, Y.; Du, Y.; Zhou, K.; Wang, J.; Zhao, W. X.; and Wen, J.-R. 2023. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*.
- Liu, H.; Li, C.; Li, Y.; Li, B.; Zhang, Y.; Shen, S.; and Lee, Y. J. 2024a. LLaVA-NeXT: Improved reasoning, OCR, and world knowledge.
- Liu, X.; Li, Z.; He, Z.; Li, P.; Xia, S.; Cui, X.; Huang, H.; Yang, X.; and He, R. 2025. Video-SafetyBench: A Benchmark for Safety Evaluation of Video LLMs. *arXiv preprint arXiv:2505.11842*.
- Liu, X.; Zhu, Y.; Gu, J.; Lan, Y.; Yang, C.; and Qiao, Y. 2024b. Mm-safetybench: A benchmark for safety evaluation of multimodal large language models. In *European Conference on Computer Vision*, 386–403. Springer.
- Liu, Y.; Li, S.; Liu, Y.; Wang, Y.; Ren, S.; Li, L.; Chen, S.; Sun, X.; and Hou, L. 2024c. TempCompass: Do Video LLMs Really Understand Videos? *arXiv preprint arXiv:2403.00476*.
- Mazeika, M.; Phan, L.; Yin, X.; Zou, A.; Wang, Z.; Mu, N.; Sakhaee, E.; Li, N.; Basart, S.; Li, B.; et al. 2024. Harmbench: A standardized evaluation framework for automated red teaming and robust refusal. *arXiv preprint arXiv:2402.04249*.

- Miao, Y.; Zhu, Y.; Yu, L.; Zhu, J.; Gao, X.-S.; and Dong, Y. 2024. T2vsafetybench: Evaluating the safety of text-to-video generative models. *Advances in Neural Information Processing Systems*, 37: 63858–63872.
- Mireshg, N.; Kim, H.; Zhou, X.; Tsvetkov, Y.; Sap, M.; Shokri, R.; and Choi, Y. 2023. Can llms keep a secret? testing privacy implications of language models via contextual integrity theory. *arXiv preprint arXiv:2310.17884*.
- Nadeem, M.; Bethke, A.; and Reddy, S. 2020. StereoSet: Measuring stereotypical bias in pretrained language models. *arXiv preprint arXiv:2004.09456*.
- Nan, K.; Xie, R.; Zhou, P.; Fan, T.; Yang, Z.; Chen, Z.; Li, X.; Yang, J.; and Tai, Y. 2024. OpenVid-1M: A Large-Scale High-Quality Dataset for Text-to-video Generation. *arXiv preprint arXiv:2407.02371*.
- Nangia, N.; Vania, C.; Bhalerao, R.; and Bowman, S. R. 2020. CrowS-pairs: A challenge dataset for measuring social biases in masked language models. *arXiv preprint arXiv:2010.00133*.
- Ning, M.; Zhu, B.; Xie, Y.; Lin, B.; Cui, J.; Yuan, L.; Chen, D.; and Yuan, L. 2023. Video-bench: A comprehensive benchmark and toolkit for evaluating video-based large language models. *arXiv preprint arXiv:2311.16103*.
- Park, J.; Jang, K. J.; Alasaly, B.; Mopidevi, S.; Zolensky, A.; Eaton, E.; Lee, I.; and Johnson, K. 2025. Assessing modality bias in video question answering benchmarks with multimodal large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 19821–19829.
- Saravanan, D.; Singh, D.; Gupta, V.; Khan, Z.; Gandhi, V.; and Tapaswi, M. 2024. VELOCITI: Can Video-Language Models Bind Semantic Concepts through Time? *arXiv preprint arXiv:2406.10889*.
- Sharma, T.; Stangl, A.; Zhang, L.; Tseng, Y.-Y.; Xu, I.; Findlater, L.; Gurari, D.; and Wang, Y. 2023. Disability-first design and creation of a dataset showing private visual information collected with people who are blind. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, 1–15.
- Socher, R.; Perelygin, A.; Wu, J.; Chuang, J.; Manning, C. D.; Ng, A.; and Potts, C. 2013. Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, 1631–1642. Seattle, Washington, USA: Association for Computational Linguistics.
- Sun, L.; Huang, Y.; Wang, H.; Wu, S.; Zhang, Q.; Gao, C.; Huang, Y.; Lyu, W.; Zhang, Y.; Li, X.; et al. 2024. Trustllm: Trustworthiness in large language models. *arXiv preprint arXiv:2401.05561*, 3.
- Szegedy, C.; Zaremba, W.; Sutskever, I.; Bruna, J.; Erhan, D.; Goodfellow, I.; and Fergus, R. 2013. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*.
- Tang, Y.; Bi, J.; Xu, S.; Song, L.; Liang, S.; Wang, T.; Zhang, D.; An, J.; Lin, J.; Zhu, R.; et al. 2025. Video understanding with large language models: A survey. *IEEE Transactions on Circuits and Systems for Video Technology*.
- Team, G.; Georgiev, P.; Lei, V. I.; Burnell, R.; Bai, L.; Gulati, A.; Tanzer, G.; Vincent, D.; Pan, Z.; Wang, S.; et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Wang, W.; He, Z.; Hong, W.; Cheng, Y.; Zhang, X.; Qi, J.; Gu, X.; Huang, S.; Xu, B.; Dong, Y.; et al. 2024a. Lvbench: An extreme long video understanding benchmark. *arXiv preprint arXiv:2406.08035*.
- Wang, Y.; Hu, W.; Dong, Y.; Liu, J.; Zhang, H.; and Hong, R. 2025a. Align is not enough: Multimodal universal jailbreak attack against multimodal large language models. *IEEE Transactions on Circuits and Systems for Video Technology*.
- Wang, Y.; Song, J.; Gao, Y.; Wang, X.; Yao, Y.; Teng, Y.; Ma, X.; Wang, Y.; and Jiang, Y.-G. 2025b. SafeVid: Toward Safety Aligned Video Large Multimodal Models. *arXiv preprint arXiv:2505.11926*.
- Wang, Y.; Wang, Y.; Zhao, D.; Xie, C.; and Zheng, Z. 2024b. Videohalluciner: Evaluating intrinsic and extrinsic hallucinations in large video-language models. *arXiv preprint arXiv:2406.16338*.
- Ying, Z.; Liu, A.; Liang, S.; Huang, L.; Guo, J.; Zhou, W.; Liu, X.; and Tao, D. 2024. Safebench: A safety evaluation framework for multimodal large language models. *arXiv preprint arXiv:2410.18927*.
- Zhang, Y.; Huang, Y.; Sun, Y.; Liu, C.; Zhao, Z.; Fang, Z.; Wang, Y.; Chen, H.; Yang, X.; Wei, X.; et al. 2024. Benchmarking trustworthiness of multimodal large language models: A comprehensive study. *arXiv preprint arXiv:2406.07057*.
- Zhang, Y.; Huang, Y.; Wang, Y.; Sun, Y.; Liu, C.; Zhao, Z.; Fang, Z.; Chen, H.; Yang, X.; Wei, X.; et al. 2025. Unveiling trust in multimodal large language models: Evaluation, analysis, and mitigation. *arXiv preprint arXiv:2508.15370*.
- Zhou, L.; Xu, C.; and Corso, J. J. 2018. Towards Automatic Learning of Procedures From Web Instructional Videos. In *AAAI Conference on Artificial Intelligence*, 7590–7598.