

VisualActBench: Can VLMs See and Act like a Human?

Daoan Zhang^{✉*} Pai Liu^{✉*} Xiaofei Zhou^{✉*} Yuan Ge^{ℳ*} Guangchen Lan^{ℙ*}
Jing Bi[✉] Christopher Brinton^ℙ Ehsan Hoque[✉] Jiebo Luo[✉]

[✉]University of Rochester ^ℙPurdue University ^ℳNortheastern University

dayu@udel.edu

{antoinesimoulin,xiaoyiliu,yuweicao,zhaoputeng,xinqian,glyang}@meta.com

{til040,jmcauley}@ucsd.edu {daoan.zhang,jluo}@rochester.edu

Abstract

Vision-Language Models (VLMs) have achieved impressive progress in perceiving and describing visual environments. However, their ability to proactively reason and act based solely on visual inputs, without explicit textual prompts, remains underexplored. We introduce a new task, **Visual Action Reasoning**, and propose **VisualActBench**, a large-scale benchmark comprising 1,074 videos and 3,733 human-annotated actions across four real-world scenarios. Each action is labeled with an Action Prioritization Level (APL) and a proactive/reactive type to assess models' human-aligned reasoning and value sensitivity. We evaluate 29 VLMs on VisualActBench and find that while frontier models like GPT-4o demonstrate relatively strong performance, a significant gap remains compared to human-level reasoning—particularly in generating proactive, high-priority actions. Our results highlight limitations in current VLMs' ability to interpret complex context, anticipate outcomes, and align with human decision-making frameworks. VisualActBench establishes a comprehensive foundation for assessing and improving the real-world readiness of proactive, vision-centric AI agents.

1 Introduction

The emergence of Vision-Language Models (VLMs) (Liu et al., 2023, 2024a; Bai et al., 2023; Zhu et al., 2023) has greatly improved their capability to perceive and comprehend the open world. However, the advent of AGI demands a paradigm shift for VLMs: these systems must evolve from reactive language-mediated agents to vision-centric proactive entities—that is, from reliance on intermediate language representations to interpret, respond to, and act upon visual inputs, to being capable of autonomously sensing environmental dynamics through visual streams and initiating

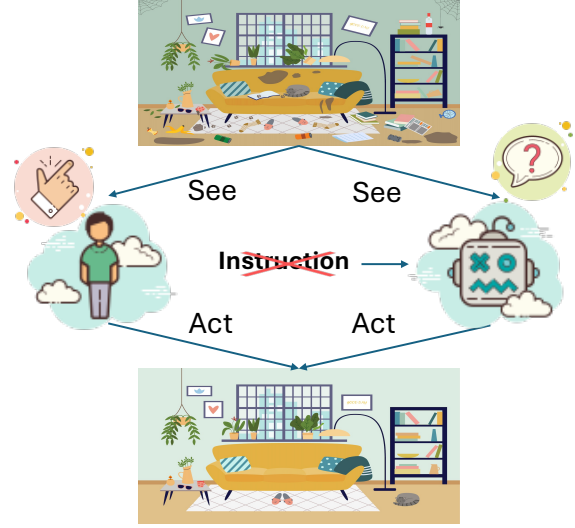


Figure 1: For humans, seeing a cluttered room naturally prompts the intention to tidy it up. However, in the absence of explicit instructions such as "tidy up the room," VLMs, relying solely on their own capabilities, infer the relevant action based on visual cues.

context-aware actions without linguistic intermediation. Comparatively, the latter approach aligns more closely with the way humans predominantly process information—through vision rather than language. For instance, humans can avoid obstacles while walking or catch a ball in mid-air by relying solely on visual perception and prior knowledge, without the need for linguistic reasoning. This vision-centered processing enables rapid decision-making and instantaneous responses, which are crucial for applications such as robotics, autonomous vehicles, and surveillance systems. However, the proactive performance of current VLMs under purely visual input conditions remains largely unknown, particularly regarding whether VLMs can generate human-consistent responses when presented with visual-only information.

Therefore, the first research question we aim to address in this study is: **Can VLMs emulate**

*Equal contribution.

human decision-making in open-world contexts with accurate and context-aware instructions?

Another significant aspect of human decision-making lies in the role of behavioral value systems (Schwartz, 2013). Unlike purely reactive responses, human actions are often guided by internalized values, ethical considerations, and context-specific objectives, which influence choices and prioritize outcomes beyond immediate sensory input. For instance, when crossing a busy street, a pedestrian not only processes visual cues like traffic lights and oncoming vehicles but also integrates learned behavioral norms—such as waiting for a green signal or yielding to vulnerable individuals. This value-based decision-making enables humans to navigate complex social and moral environments, ensuring their actions align with long-term goals, safety, and social acceptability. In contrast, current VLMs primarily focus on understanding and describing visual scenes without inherently embedding or prioritizing actions based on human-aligned value systems. The models are trained to optimize performance on prompting tasks, such as object detection, image captioning, or instruction generation, often without accounting for the ethical or social implications of the responses they produce. This limitation raises an important question regarding their deployment in real-world settings, particularly in scenarios requiring nuanced decision-making, safety considerations, and adherence to social norms. Thus, the second research question this study seeks to explore is: **Do VLMs exhibit behavior consistent with general human value systems when generating instructions in open-world contexts?**

To address the two key challenges outlined earlier—evaluating VLMs based solely on visual inputs and assessing their alignment with human preferences—we introduce a novel task, **Visual Action Reasoning**, along with its corresponding benchmark, **VisualActBench**. This task departs from prior work by being strictly *visual-centric*, meaning that all the information required for the model to generate an appropriate response is embedded solely within the visual input, without relying on any explicit textual instructions or prompts. The goal is to assess whether VLMs can extract meaningful cues from visual scenes and reason about the most suitable actions in a manner that reflects human-like understanding and decision-making.

To construct the benchmark, we define four representative real-world scenarios: *Dynamic Naviga-*

Table 1: A comparative overview of various benchmarks across several dimensions, such as data format (image **I** or video **V**), the size of dataset (**#Data**), Visual-centric(indicated by **VC.**), Action Reasoning Generation(indicated by **ARG.**), Human Alignment(Action prioritization level, indicated by **HA.**), Open-end question(indicated by **OQ.**) and the annotation method (manual or automatic/manual, indicated by **Anno.**).

Benchmark	I/V	#Data	VC.	ARG.	HA.	OQ.	Anno.
Image Datasets (I)							
Winoground (Thrush et al., 2022)	I	400	×	×	×	×	M
MME (Fu et al., 2023)	I	1.1k	×	×	×	✓	M
MMBench (Liu et al., 2025a)	I	1.7k	×	×	×	×	A+M
MMComposition (Hua et al., 2024)	I	4.3k	×	×	×	×	M
Video Datasets (V)							
MSRVTT (Xu et al., 2016)	V	10k	✓	×	×	✓	A+M
MSVD-QA (Xu et al., 2017)	V	504	×	×	×	×	A
MSRVTT-QA (Xu et al., 2017)	V	2.9k	×	×	×	×	A
TGIF-QA (Jang et al., 2017)	V	9.6k	×	×	×	×	A
TVQA (Lei et al., 2018)	V	2.2k	×	×	×	×	A+M
ActivityNet-QA (Yu et al., 2019)	V	5.8k	×	×	×	×	M
NExT-QA (Xiao et al., 2021)	V	1k	×	×	×	×	A
AutoEval-Video (Chen et al., 2023)	V	327	×	×	×	×	A+M
Video-Bench (Ning et al., 2023)	V	5.9k	×	×	×	×	A+M
LVBench (Wang et al., 2024c)	V	500	×	×	×	×	M
MVBench (Li et al., 2024c)	V	3.6k	×	×	×	✓	A+M
MovieChat-1k (Song et al., 2024)	V	100	×	×	×	✓	M
TempCompass (Liu et al., 2024b)	V	410	×	×	×	×	M
Video-MME (Fu et al., 2024)	V	900	×	×	×	×	M
HumanActBench	V	1074	✓	✓	✓	✓	A+M

tion, Home Service, Safety and Monitoring, and Human-Machine Interaction. We sample 1,074 videos from existing datasets covering these contexts and manually annotate a total of 3,733 actions. Each action is paired with an **Action Prioritization Level (APL)**—a metric designed to reflect human-aligned preference levels based on contextually appropriate responses. Additionally, actions are categorized into either *proactive* or *reactive* types, providing a lens to assess the model’s tendency toward initiative-taking behavior. This allows us to measure not only the correctness of action generation but also the extent to which models exhibit *human-aligned proactiveness*.

As shown in Table 1, VisualActBench differs from traditional video benchmarks such as captioning or question-answering datasets by emphasizing action reasoning generation (ARG) and preference alignment. Instead of prompting models with textual questions or commands, our setting requires them to autonomously analyze the scene, anticipate contextually desirable outcomes, and act accordingly. This design makes VisualActBench a comprehensive and challenging benchmark that holistically evaluates both the perceptual and cognitive aspects of VLM behavior in human-like tasks.

We evaluate 29 VLMs on VisualActBench, including 24 open-source models and 5 proprietary models. Our findings reveal that even the best-

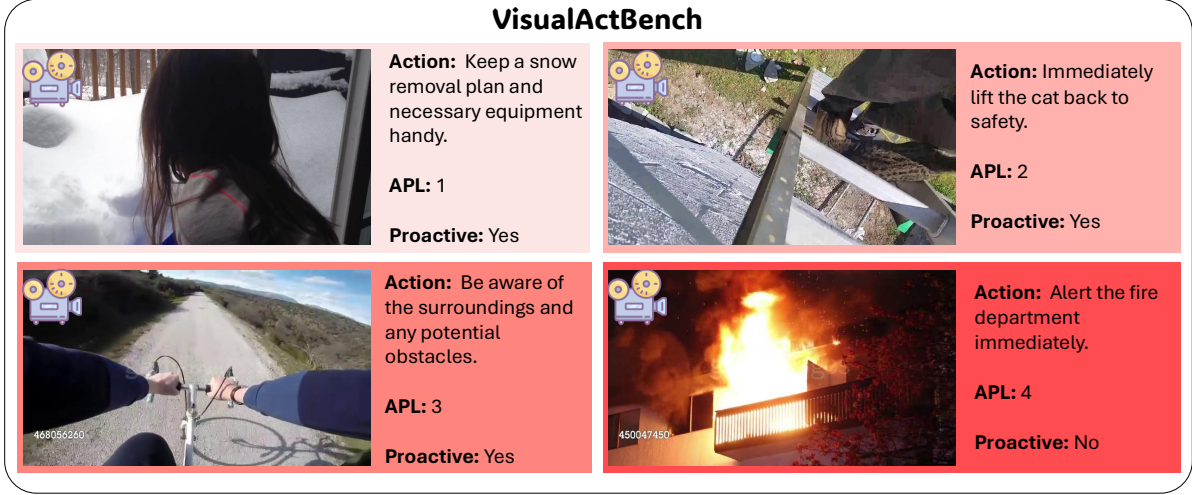


Figure 2: Examples from **VisualActBench**, showcasing diverse real-world scenarios and the corresponding proactive actions with varying Action Priority Levels (APL). Each frame includes a proposed action, its associated APL, and whether the action is considered proactive.

performing VLMs exhibit a substantial gap compared to human performance on this task. Moreover, since the Visual Action Reasoning task requires both visual perception and LLM-driven reasoning capabilities, it serves as a more comprehensive evaluation of VLMs’ overall ability.

Our main contributions are:

- **A new vision-centric task:** We propose *Visual Action Reasoning*, which challenges VLMs to generate context-aware, proactive actions directly from visual input—without any textual prompts—emulating human decision-making in dynamic environments.
- **A human-aligned benchmark:** We introduce *VisualActBench*, covering four real-world scenarios with over 3,700 annotated actions, each labeled by Action Prioritization Level (APL) and proactive/reactive type to assess value alignment and behavioral quality.
- **A comprehensive model evaluation:** We benchmark 30 VLMs and reveal significant gaps in proactiveness, value alignment, and abstract reasoning, offering critical insights into the limitations of current models and guiding future development toward real-world deployment.

2 VisualActBench

2.1 Video Collection

We sampled the videos from 4 datasets, Kinetics (Kay et al., 2017), Moments in Time (Monfort

et al., 2019), DADA-2000 (Fang et al., 2019) and UCF-Crime (Sultani et al., 2018) datasets. Notably, we select video clips that are not too long and, as much as possible, contain only a single scene. This approach helps avoid an excessive number of actions and complications related to action sequencing that may arise from multiple scenes.

2.2 Benchmark Statistics

Figure 3 provides a quantitative overview of the composition of **VisualActBench**. As shown in the leftmost pie chart, the dataset contains a balanced distribution across four key scenarios: Dynamic Navigation (364 videos), Home Service (454), Safety and Monitoring (126), and Human-Machine Interaction (130), ensuring diverse contextual coverage. The adjacent chart shows the Action Prioritization Level (APL) distribution, with the majority of videos concentrated in Level 0 (455) and Level 2 (434), while higher-risk or more critical scenarios (Level 3 and 4) are less common, reflecting the natural rarity of high-priority decision-making moments in everyday environments.

On the action side, the benchmark contains a total of 3,733 actions, as depicted in the two right-side pie charts. Notably, proactive actions significantly outnumber reactive ones (2,685 vs. 1,048), indicating that the benchmark emphasizes initiative-driven reasoning over passive response. The APL distribution of actions shows a smooth decline from Level 0 (1,675) to Level 4 (69), further confirming that while lower-priority actions are more frequent, higher-priority and more critical

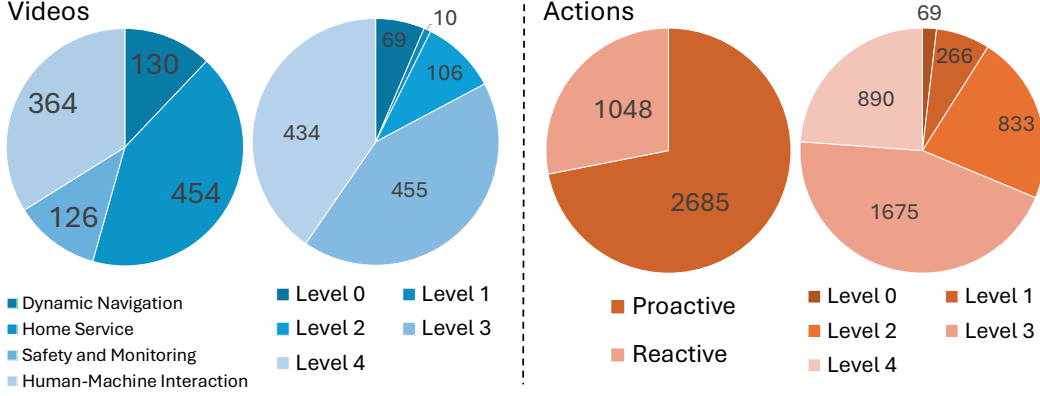


Figure 3: Distribution of videos and actions in **VisualActBench**. The left two charts show the number of videos categorized by scenario type (Dynamic Navigation, Home Service, Safety and Monitoring, and Human-Machine Interaction) and by Action Priority Level (APL 0–4). The right two charts illustrate the number of actions categorized by proactivity (Proactive vs. Reactive) and by APL.

actions are still well-represented. This balanced distribution supports comprehensive evaluation of a model’s ability to reason across a spectrum of urgency and proactiveness.

2.3 Evaluation Metrics

To evaluate model performance in VisualActBench, we employ a set of metrics that capture both textual alignment and human-aligned reasoning quality. Specifically, we calculate standard precision, recall, and F1-score based on action matching, along with two additional metrics: weighted alignment score and scale score. All metrics are reported at the average level across four scenario classes.

Exact Matching-based Metrics. We use cosine similarity over sentence embeddings (produced by a SentenceTransformer model) to compute the similarity between ground-truth and predicted actions. For a pair of predicted action p_i and ground-truth action g_j , we compute a similarity matrix S where

$$S_{i,j} = \text{cosine_similarity}(p_i, g_j).$$

Actions are matched using the Hungarian algorithm with a similarity threshold τ (e.g., $\tau = 0.5$). Based on the matched pairs:

- True Positives (TP): matched action pairs above the threshold.
- False Positives (FP): unmatched predicted actions.
- False Negatives (FN): unmatched ground-truth actions.

From these, we compute Precision, Recall and F1:

Weighted Matching Metrics. To incorporate human value alignment into the evaluation, we propose a weighted F1 variant that considers how well the predicted action reflects the importance level (scale) of the ground-truth action. For each matched pair (p_i, g_j) with associated scale levels s_i and s_j , we define a match weight as:

$$w_{i,j} = 1 - \frac{|s_i - s_j|}{4},$$

where the scale difference is normalized by the maximum level (4). The weighted precision and recall are then:

$$\text{Weighted Precision} = \frac{\sum_{(i,j) \in \text{matches}} w_{i,j}}{|\text{Predicted}|}, \quad (1)$$

$$\text{Weighted Recall} = \frac{\sum_{(i,j) \in \text{matches}} w_{i,j}}{|\text{Ground Truth}|}. \quad (2)$$

and the weighted F1 score follows similarly.

Average Scale Score. To further evaluate human alignment, we introduce the average scale score, measuring the closeness of predicted action scale to ground-truth:

$$\text{Scale Score} = \frac{1}{|\text{matches}|} \sum_{(i,j) \in \text{matches}} \left(1 - \frac{|s_i - s_j|}{4} \right).$$

This metric penalizes mismatches in action priority even if the actions are textually similar, reinforcing the importance of value-sensitive reasoning.

For fairness across diverse scenarios, we report averaged metrics by computing the mean of each

metric across the four predefined classes, which ensures equal weight for each scenario type regardless of sample size.

3 Experiment Results

We evaluate over 30 state-of-the-art vision-language models (VLMs) on **VisualActBench** to assess their ability to generate context-aware actions across diverse real-world scenarios. This section presents key quantitative results and insights, highlighting both overall performance and scenario-specific trends that reveal current strengths and limitations in multimodal reasoning.

3.1 Weighted Matching Metric Comparison

Table 2 presents a comprehensive quantitative evaluation of 30 vision-language models (VLMs) on **VisualActBench**, focusing on their ability to generate contextually appropriate actions across four real-world scenarios—Dynamic Navigation, Home Service, Safety and Monitoring, and Human-Machine Interaction. For each model, we report precision, recall, and F1 scores scaled to [0, 100], with the top three scores per column highlighted for clarity. Among all models, **GPT-4o** achieves the highest overall F1 score of 36.7, significantly outperforming both its smaller variant GPT-4o-mini (32.7) and other strong baselines such as InternVL2.5-38B (33.7) and Aria (32.5). GPT-4o consistently ranks in the top three across nearly all categories, confirming its advanced ability in not just visual perception but also high-level reasoning and decision-making. Interestingly, GPT-4o-mini—despite its compact size—still achieves competitive recall, particularly in Home Service and Human-Machine Interaction, indicating a strong balance between efficiency and reasoning ability.

We also observe a stark performance gap between the top-performing models and the rest of the field. Many popular open-source VLMs, including mPLUG-Owl, VideoLLaMA, and various versions of InternVL below 8B, report overall F1 scores below 20, revealing fundamental limitations in generating high-priority or proactive actions. These models often exhibit high precision but poor recall, suggesting a tendency to output conservative or generic actions, thereby missing more critical or context-sensitive responses. This imbalance indicates a failure in identifying nuanced cues that signal the need for proactive intervention, such as anticipating a fall, preventing a fire, or offering

assistance.

Scenario-specific analysis further reveals that models perform best in *Dynamic Navigation* and *Safety and Monitoring*, where visual cues are often explicit and well-structured (e.g., obstacles, road layouts, fire). In contrast, the *Home Service* and *Human-Machine Interaction* domains yield lower scores across all models. These scenarios typically require the model to understand abstract intent, anticipate user needs, or consider long-term outcomes—capabilities that current VLMs largely lack. For example, generating actions like “adjust the lighting for comfort” or “check if the device is overheating” demands both commonsense and forward-looking reasoning.

Moreover, the introduction of the **Action Prioritization Level (APL)** and proactive/reactive classification in VisualActBench enables us to go beyond surface-level accuracy and examine alignment with human values. Our results show that many models default to low-APL or reactive responses, failing to reflect human preferences that favor preemptive and high-impact interventions. This aligns with our earlier findings that VLMs tend to act as passive observers rather than active decision-makers. The challenge is not only in recognizing what is happening in the scene, but in inferring what *ought to be done*—a fundamentally different and cognitively demanding capability.

Qualitative inspection reveals typical failure modes such as overly generic actions (“observe the scene”), redundant outputs (“check again”), or misaligned priorities (e.g., focusing on minor details over safety-critical cues). These errors suggest that VLMs often rely on shallow correlations rather than genuine situational understanding.

3.2 Average Scale Score Evaluation

Table 3 shows the normalized scale scores of different VLMs across four core scenarios, offering a more human-aligned perspective on model performance. While **GPT-4o** clearly leads with the highest overall score (66.4), models like **GPT-4o-mini**, **Aria**, and **InternVL2.5-38B** also achieve comparably strong performance, suggesting that multiple models can produce contextually appropriate actions even if they differ in architecture and size. Notably, the Safety and Monitoring (SM) scenario sees generally higher scores, likely due to its more explicit visual cues. In contrast, Human-Machine Interaction (HMI) remains challenging across the board, reflecting the difficulty of interpreting in-

Table 2: Comprehensive evaluation results (scores scaled to 100 and rounded to one decimal), sorted by Overall F1 (descending). In each column the highest value (top 1) is highlighted in gold, the second highest (top 2) in silver and the third highest (top 3) in bronze.

Model	Dynamic Navigation			Home Service			Safety and Monitoring			Human-Mach. Inter.			Overall		
	Prec.	Recall	F1	Prec.	Recall	F1	Prec.	Recall	F1	Prec.	Recall	F1	Prec.	Recall	F1
GPT-4o (OpenAI, 2024)	32.6	40.9	34.6	38.5	41.5	38.4	34.5	40.9	35.8	37.5	42.6	38.2	35.8	41.5	36.7
InternVL2.5-38B (Chen et al., 2025)	40.7	33.9	33.9	41.4	28.0	30.9	44.1	32.3	34.6	45.7	32.7	35.3	43.0	31.7	33.7
GPT4o-mini (OpenAI, 2024)	25.4	33.1	27.4	37.5	42.2	38.0	27.6	34.1	29.2	35.1	41.0	36.1	31.4	37.6	32.7
Aria (Li et al., 2024b)	35.7	36.2	32.9	35.0	30.0	30.3	57.2	38.2	34.7	36.6	34.6	32.8	36.4	33.8	32.5
Gemini-1.5-Flash-8B (DPM, 2024)	32.6	40.3	33.7	28.4	28.9	26.1	32.5	35.4	31.1	32.6	36.1	31.3	31.5	35.2	30.6
LLaVA-OV-72B (Li et al., 2024a)	27.2	33.4	28.3	31.0	31.9	29.9	26.5	29.1	26.4	29.2	32.3	28.7	28.5	31.7	28.3
QwenVL-72B (Wang et al., 2024b)	29.1	27.2	26.6	27.9	20.9	22.6	29.9	25.1	25.9	29.4	24.1	24.9	29.1	24.3	25.0
InternVL2.5-4B (Chen et al., 2024)	44.7	22.7	27.4	40.2	16.7	22.0	39.5	18.1	23.7	39.8	19.3	23.9	41.0	19.2	24.3
Gemini-1.5-Flash (GDM, 2024a)	25.4	31.1	25.9	21.6	21.9	19.6	23.1	29.7	24.4	27.2	27.8	24.4	24.3	27.6	23.6
QwenVL-3B (Wang et al., 2024b)	39.5	21.3	24.2	27.4	13.4	15.4	31.3	16.5	18.8	35.1	18.3	20.9	33.3	17.3	19.8
NVILA-8B (Liu et al., 2025b)	24.9	22.4	21.7	21.9	16.4	17.2	23.8	18.4	19.0	21.4	18.1	18.0	23.0	18.8	19.0
Gemini-2-Flash (GDM, 2024b)	31.6	20.7	22.9	28.0	15.3	18.0	28.8	19.8	21.5	29.7	17.3	20.5	39.2	29.5	18.3
InternVL2.5-8B-MPO	30.0	16.6	19.1	27.3	13.4	16.2	32.9	15.2	19.3	27.8	13.8	16.6	29.5	14.7	17.8
MiniCPM-o (Yu et al., 2024)	30.9	20.0	22.0	25.4	14.0	16.8	29.0	16.6	20.1	28.8	16.2	18.9	38.9	28.5	16.7
InternVL2.5-8B (Chen et al., 2025)	21.9	15.5	15.9	24.1	13.1	15.4	29.3	15.3	18.8	28.2	16.3	18.3	37.4	25.9	15.0
InternVL2.5-78B(Chen et al., 2025)	15.4	15.1	14.4	18.8	13.7	15.0	14.9	11.1	12.3	20.4	16.9	17.5	17.4	14.2	14.8
VideoLLaMA2-7B (Cheng et al., 2024)	18.4	18.1	17.4	14.2	12.4	12.6	15.9	15.9	15.1	13.6	12.1	12.2	15.5	14.7	14.3
VideoLLaMA2-7B-16F	15.0	16.2	14.9	14.7	13.0	13.0	14.8	13.4	13.6	13.5	12.5	12.2	14.5	13.7	13.4
QwenVL-7B (Wang et al., 2024a)	32.0	17.7	20.9	24.1	11.4	13.9	29.1	9.7	12.8	26.1	13.5	16.1	29.4	26.3	13.1
NVILA-15B (Liu et al., 2025b)	10.5	12.5	11.2	8.8	10.7	9.0	7.9	9.9	8.3	12.4	14.2	12.7	17.2	9.9	11.8
mPLUG-Owl3-7B	21.5	11.0	13.1	14.1	5.3	7.1	18.8	8.4	10.6	20.3	9.4	11.7	18.7	8.5	10.6
MiniCPM-v (Yao et al., 2024)	15.7	9.0	10.7	12.0	6.8	8.2	16.5	8.8	10.9	10.8	16.2	10.8	20.5	14.7	8.6
InternVL2.5-2B (Chen et al., 2025)	8.8	7.3	7.0	9.8	4.7	5.5	9.2	6.1	6.7	10.1	6.5	7.0	15.5	9.5	6.1
LLaVA-OV-7B (Li et al., 2024a)	7.7	7.1	7.3	6.0	4.4	4.8	3.9	2.7	2.9	9.8	9.5	9.4	6.9	5.9	6.1
VILA-1.5-40B (Lin et al., 2024)	11.8	5.6	6.6	10.7	4.4	5.6	11.6	3.8	5.6	10.2	3.9	5.2	10.9	11.1	4.4
VideoChat-Flash-2B (Li et al., 2025)	3.3	1.5	1.9	6.9	3.1	4.0	3.3	1.8	2.2	7.0	3.1	4.0	5.1	2.4	3.0
InternVL2.5-1B (Chen et al., 2025)	1.4	0.5	0.7	3.2	1.7	2.0	3.7	2.8	2.9	2.7	2.7	2.5	2.8	1.9	2.0
VideoChat-Flash-7B (Li et al., 2025)	3.1	0.7	1.1	3.2	1.3	1.7	1.7	1.3	1.4	2.3	0.7	1.0	2.6	1.0	1.3
LLaVA-OV-0.5B (Li et al., 2024a)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.1	0.1	0.2	0.0	0.0	0.1

tent and social context. Some models show wide variance between scenarios, indicating specialization, while others maintain consistent but modest scores. Overall, the scale score metric highlights not only top-line performance but also qualitative stability and contextual sensitivity—key factors for real-world deployment.

3.3 Proactiveness of VLMs

Figure 4 reports the normalized proactive ratio for each model, indicating how well each VLM completes the proactive actions expected in the evaluation set. Rather than simply preferring proactive outputs, this metric reflects a model’s ability to detect situations that require initiative and respond accordingly.

Results show that only a handful of models—such as **InternVL2.5-8B-MPO**, **Gemini-1.5-Flash**, and **QwenVL-3B**—achieve strong proactive coverage, successfully completing over 70

These findings highlight a major gap: most VLMs still struggle to recognize and act on sit-

uations where proactive, high-priority behavior is needed—limiting their applicability in safety-critical or assistance-driven environments.

3.4 Impact of Frame Count and Model Size on Action Quality

Frame Count Analysis. The results are shown in Table. 4. The number of input frames has a notable effect on model performance, but not always positively. For LLaVA-OV, increasing frame count from 2 to 16 leads to a consistent decline in F1 and APL scores. This trend can be attributed to the nature of our benchmark videos, which typically depict a single, well-defined event. Including more frames often introduces redundant information before or after the core event, thereby diluting the model’s focus and impairing its ability to generate precise and prioritized actions. In contrast, VideoLLaMA2 demonstrates more stable performance across 8 and 16 frames, suggesting better temporal robustness. Nonetheless, the marginal differences indicate that once the core event is captured, addi-

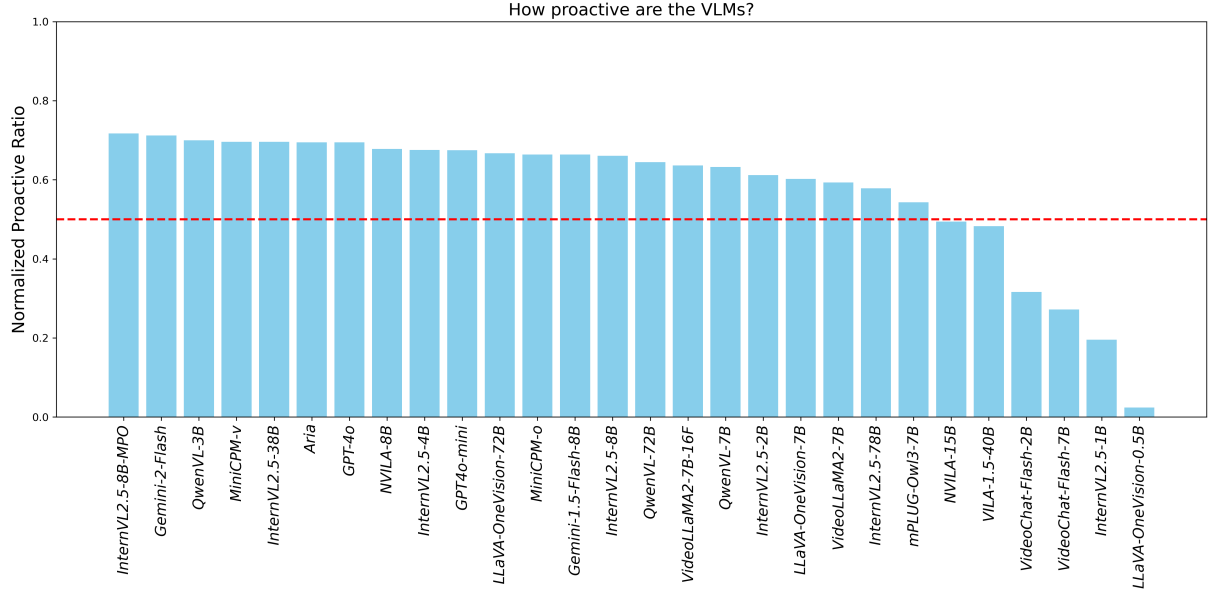


Figure 4: Normalized proactive ratios of various VLMs, reflecting their overall inclination to generate proactive rather than reactive actions across diverse scenarios. The red dashed line denotes the average proactive ratio across all evaluated models, serving as a reference for assessing the relative proactiveness and behavioral consistency of each system.

tional frames may offer limited benefit—or even harm performance if not properly handled.

Model Size Analysis. The results are shown in Table 5. Across all model families—LLaVA-OV, QwenVL, and InternVL—performance consistently improves with larger model size, both in terms of F1 score and APL. For instance, LLaVA-OV scales from near-zero performance at 0.5B to a strong F1 of 28.3 and APL of 54.0 at 72B. Similarly, InternVL shows a significant jump from 1B (F1: 2.0, APL: 3.6) to 38B (F1: 33.7, APL: 57.9).

However, performance gains are not always linear. InternVL-78B, for example, sees a drop in F1 and APL compared to InternVL-38B, suggesting diminishing returns or possible overfitting/misalignment at extreme scales. This hints that beyond a certain size, architectural design and training data quality may play a larger role than raw parameter count.

These findings collectively suggest that model performance in Visual Action Reasoning is constrained more by temporal abstraction and contextual grounding than by sheer model capacity. Future designs should thus focus on adaptive temporal encoding and selective frame reasoning, enabling models to extract causally relevant cues without redundancy.

3.5 Effect of Reinforcement Learning

As shown in 6, reinforcement learning (via Mixed Preference Optimization (MPO)) generally improves model performance, especially in larger models. For example, InternVL-38B shows notable gains in F1 (33.7 → 40.7) and APL (57.9 → 58.2), indicating better action quality and prioritization. While smaller models like InternVL-8B see mixed results, the trend suggests that RL enhances proactiveness and recall, helping models generate more complete and human-aligned actions.

4 Conclusion

In this work, we presented **VisualActBench**, a novel benchmark designed to evaluate the proactive reasoning capabilities of Vision-Language Models (VLMs) in a vision-only setting. By introducing the *Visual Action Reasoning* task, we challenge models to move beyond passive visual understanding and toward human-like, initiative-driven decision-making grounded in visual context alone. Through the annotation of 3,733 actions across four diverse scenarios, along with a novel Action Prioritization Level (APL) metric, VisualActBench provides a rich framework for evaluating both the correctness and value-alignment of model outputs. Our comprehensive evaluation of 29 state-of-the-art VLMs reveals key insights: current models struggle with proactiveness, value-sensitive prioritization, and ab-

Table 3: Scale scores (on a 0–100 scale, rounded to one decimal) for various models, sorted by Overall (from high to low). In each column, the highest value (top 1) is highlighted in gold, the second highest (top 2) in silver, and the third highest (top 3) in bronze.

Model	DN.	HS.	SM.	HML.	Overall
GPT-4o	69.4	64.0	70.4	61.9	66.4
GPT4o-mini	57.1	64.0	60.0	60.4	60.4
Aria	61.6	57.2	62.2	59.7	60.2
IntVL2.5-38B	58.1	55.7	58.3	59.6	57.9
Gemini-1.5-Flash-8B	59.6	53.8	59.0	58.7	57.8
LLaVA-OneVision-72B	53.2	56.8	52.2	53.7	54.0
Gemini-1.5-Flash	51.0	45.7	50.6	49.9	49.3
InternVL2.5-4B	52.4	46.2	45.6	46.3	47.6
QwenVL-72B	47.7	42.6	42.8	44.5	44.4
QwenVL-3B	49.4	34.0	38.4	40.8	40.6
NVILA-8B	47.9	37.0	38.3	36.9	40.0
Gemini-2-Flash	43.4	36.1	39.8	39.2	39.6
MiniCPM-o	41.1	36.3	36.7	38.9	38.2
InternVL2.5-8B-MPO	40.1	36.1	38.8	34.0	37.2
InternVL2.5-8B	33.1	31.5	37.8	37.4	34.9
QwenVL-7B	39.1	29.1	25.8	29.4	30.9
VideoLLaMA2-7B	33.7	27.7	31.5	24.6	29.4
VideoLLaMA2-7B-16F	34.0	29.0	28.5	25.4	29.2
InternVL2.5-78B	23.3	26.4	20.1	26.1	24.0
mPLUG-Owl3-7B	25.8	15.8	20.7	22.2	21.1
MiniCPM-v	21.6	17.2	24.5	20.5	21.0
NVILA-15B	16.1	15.4	12.9	17.2	15.4
InternVL2.5-2B	14.6	13.9	14.2	15.5	14.5
VILA-1.5-40B	12.9	12.0	5.6	12.3	12.0
VideoChat-Flash-2B	4.2	6.9	2.2	4.2	9.0
LLaVA-OneVision-7B	9.2	4.4	2.9	4.8	7.8
InternVL2.5-1B	1.9	3.2	2.0	5.1	3.6
VideoChat-Flash-7B	3.1	3.2	1.4	2.2	2.6
LLaVA-OneVision-0.5B	0.0	0.0	0.0	2.1	1.5

strat reasoning. Although large-scale proprietary models such as GPT-4o perform best, even they fall short of human-level performance. Notably, increasing model size and applying reinforcement learning techniques improve action quality and prioritization, but diminishing returns are observed beyond a certain scale. Overall, **VisualActBench** exposes critical limitations in current VLMs’ ability to act autonomously and appropriately in real-world, open-ended contexts. We hope this benchmark will serve as a stepping stone for the development of more aligned, robust, and contextually aware vision-language agents capable of operating without explicit instructions—truly seeing and acting like humans.

Model	Frm.	Prec.	Recall	F1	APL
LLaVA-OV	2	18.6	8.7	10.6	17.3
	4	14.6	6.3	8.0	13.8
	8	6.9	5.9	6.1	7.8
	16	9.7	4.4	5.3	10.0
VideoLLaMA2	8	15.5	14.7	14.3	29.4
	16	14.5	13.7	13.4	29.2

Table 4: Effect of varying input frame counts on model performance, showing that increasing frames from 2 to 16 does not always improve results—particularly for LLaVA-OV, where redundant visual context often reduces precision and APL, while VideoLLaMA2 maintains more stable performance across different temporal inputs.

Model	size	Prec.	Recall	F1	APL
LLaVA-OV	0.5B	0.0	0.0	0.1	1.5
	7B	6.9	5.9	6.1	7.8
	72B	28.5	31.7	28.3	54.0
QwenVL	3B	33.3	17.3	19.8	40.6
	7B	29.4	26.3	13.1	30.9
	72B	29.1	24.3	25.0	44.4
InternVL	1B	2.8	1.9	2.0	3.6
	2B	15.5	9.5	6.1	14.5
	4B	41.0	19.2	24.3	47.6
	8B	37.4	25.9	15.0	34.9
	38B	43.0	31.7	33.7	57.9
	78B	17.4	14.2	14.8	24.0

Table 5: Comparison of model sizes across different architectures, illustrating that larger models generally yield higher F1 and APL scores, though gains taper off at extreme scales (e.g., InternVL-78B), suggesting diminishing returns and the growing importance of data quality and training strategy beyond parameter count.

Model	RL	Prec.	Recall	F1	APL
InternVL-4B	-	41.0	19.2	24.3	47.6
	MPO	49.3	20.1	26.5	48.3
InternVL-8B	-	37.4	25.9	15.0	34.9
	MPO	29.5	14.7	19.8	37.2
InternVL-38B	-	43.0	31.7	33.7	57.9
	MPO	49.0	39.9	40.7	58.2

Table 6: Impact of reinforcement learning via MPO on proactive reasoning, showing consistent improvements in both F1 and APL—particularly for large-scale models—indicating that RL enhances context awareness and alignment with human-preferred action prioritization.

References

- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*.
- Xiuyuan Chen, Yuan Lin, Yuchen Zhang, and Weiran Huang. 2023. Autoeval-video: An automatic benchmark for assessing large vision language models in open-ended video question answering. *arXiv preprint arXiv:2311.14906*.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, Lixin Gu, Xuehui Wang, Qingyun Li, Yimin Ren, Zixuan Chen, Jiapeng Luo, Jiahao Wang, Tan Jiang, Bo Wang, Conghui He, Botian Shi, Xingcheng Zhang, Han Lv, Yi Wang, Wenqi Shao, Pei Chu, Zhongying Tu, Tong He, Zhiyong Wu, Huipeng Deng, Jiaye Ge, Kai Chen, Kaipeng Zhang, Limin Wang, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhai Wang. 2025. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *Preprint*, arXiv:2412.05271.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, . . . , and Wenhai Wang. 2024. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*.
- Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, et al. 2024. VideoLLaMA 2: Advancing spatial-temporal modeling and audio understanding in video-LLMs. *arXiv preprint arXiv:2406.07476*.
- DPM. 2024. Gemini 1.5 flash-8b. Tech. announcement.
- Jianwu Fang, Dingxin Yan, Jiahuan Qiao, Jianru Xue, He Wang, and Sen Li. 2019. Dada-2000: Can driving accident be predicted by driver attention? analyzed by a benchmark. In *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pages 4303–4309. IEEE.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. 2023. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*.
- Chaoyou Fu, Yuhao Dai, Yondong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. 2024. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*.
- GDM. 2024a. Gemini 1.5 flash (lightweight multimodal model). Google I/O Announcement.
- GDM. 2024b. Introducing gemini 2.0: our new ai model for the agentic era. Google I/O Announcement.
- Hang Hua, Yunlong Tang, Ziyun Zeng, Liangliang Cao, Zhengyuan Yang, Hangfeng He, Chenliang Xu, and Jiebo Luo. 2024. Mmcomposition: Revisiting the compositionality of pre-trained vision-language models. *arXiv preprint arXiv:2410.09733*.
- Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and Gunhee Kim. 2017. Tgif-qa: Toward spatio-temporal reasoning in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2758–2766.
- Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. 2017. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*.
- Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L Berg. 2018. Tvqa: Localized, compositional video question answering. In *EMNLP*.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, et al. 2024a. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.
- Dongxu Li, Yudong Liu, Haoning Wu, Yue Wang, Zhiqi Shen, Bowen Qu, Xinyao Niu, Guoyin Wang, Bei Chen, and Junnan Li. 2024b. Aria: An open multimodal native mixture-of-experts model. *arXiv preprint arXiv:2410.05993*.
- Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. 2024c. Mvbench: A comprehensive multi-modal video understanding benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22195–22206.
- Xinhao Li, Yi Wang, Jiashuo Yu, Xiangyu Zeng, Yuhao Zhu, Haian Huang, Jianfei Gao, Kunchang Li, Yinan He, Chenting Wang, Yu Qiao, Yali Wang, and Limin Wang. 2025. Videochat-flash: Hierarchical compression for long-context video modeling. *arXiv preprint arXiv:2501.00574*.
- Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad Shoeybi, and Song Han. 2024. VILA: On pre-training for visual language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916.

- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. 2025a. Mmbench: Is your multi-modal model an all-around player? In *European Conference on Computer Vision*, pages 216–233. Springer.
- Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and Lu Hou. 2024b. Tempcompass: Do video llms really understand videos? *arXiv preprint arXiv:2403.00476*.
- Zhijian Liu, Ligeng Zhu, Baifeng Shi, Zhuoyang Zhang, Yuming Lou, Shang Yang, Haocheng Xi, Shiyi Cao, Yuxian Gu, Dacheng Li, Xiuyu Li, Yunhao Fang, Yukang Chen, Cheng-Yu Hsieh, De-An Huang, An-Chieh Cheng, Vishwesh Nath, Jinyi Hu, Sifei Liu, Ranjay Krishna, Daguang Xu, Xiaolong Wang, Pavlo Molchanov, Jan Kautz, Hongxu Yin, Song Han, and Yao Lu. 2025b. *Nvila: Efficient frontier visual language models*. *Preprint*, arXiv:2412.04468.
- Mathew Monfort, Alex Andonian, Bolei Zhou, Kandan Ramakrishnan, Sarah Adel Bargal, Tom Yan, Lisa Brown, Quanfu Fan, Dan Gutfreund, Carl Vondrick, et al. 2019. Moments in time dataset: one million videos for event understanding. *IEEE transactions on pattern analysis and machine intelligence*, 42(2):502–508.
- Munan Ning, Bin Zhu, Yujia Xie, Bin Lin, Jiayi Cui, Lu Yuan, Dongdong Chen, and Li Yuan. 2023. Video-bench: A comprehensive benchmark and toolkit for evaluating video-based large language models. *arXiv preprint arXiv:2311.16103*.
- OpenAI. 2024. *Gpt-4o mini: advancing cost-efficient intelligence*.
- OpenAI. 2024. *Gpt-4o system card*. *Preprint*, arXiv:2410.21276.
- Shalom Schwartz. 2013. Value priorities and behavior: Applying a theory of integrated value systems. In *The psychology of values*, pages 1–24. Psychology Press.
- Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe Chi, Xun Guo, Tian Ye, Yanting Zhang, et al. 2024. Moviechat: From dense token to sparse memory for long video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18221–18232.
- Waqas Sultani, Chen Chen, and Mubarak Shah. 2018. Real-world anomaly detection in surveillance videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6479–6488.
- Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. 2022. Winoground: Probing vision and language models for visio-linguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5238–5248.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024a. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Junyang Lin, Chang Zhou, Jingren Zhou, and Jinze Lin. 2024b. Qwen2-VL: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Weihaan Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Shiyu Huang, Bin Xu, Yuxiao Dong, Ming Ding, et al. 2024c. Lvbench: An extreme long video understanding benchmark. *arXiv preprint arXiv:2406.08035*.
- Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. 2021. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9777–9786.
- Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. 2017. Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 25th ACM international conference on Multimedia*, pages 1645–1653.
- Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5288–5296.
- Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. 2024. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*.
- Tianyu Yu, Haoye Zhang, Qiming Li, Qixin Xu, Yuan Yao, Da Chen, Xiaoman Lu, Ganqu Cui, Yunkai Dang, Taiwen He, Xiaocheng Feng, Jun Song, Bo Zheng, Zhiyuan Liu, Tat-Seng Chua, and Maosong Sun. 2024. *Rlaif-v: Open-source ai feedback leads to super gpt-4v trustworthiness*. *Preprint*, arXiv:2405.17220.
- Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. 2019. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 9127–9134.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.