

OVBench: How Far is Your Video-LLMs from Real-World Online Video Understanding?

Anonymous CVPR submission

Paper ID 4962

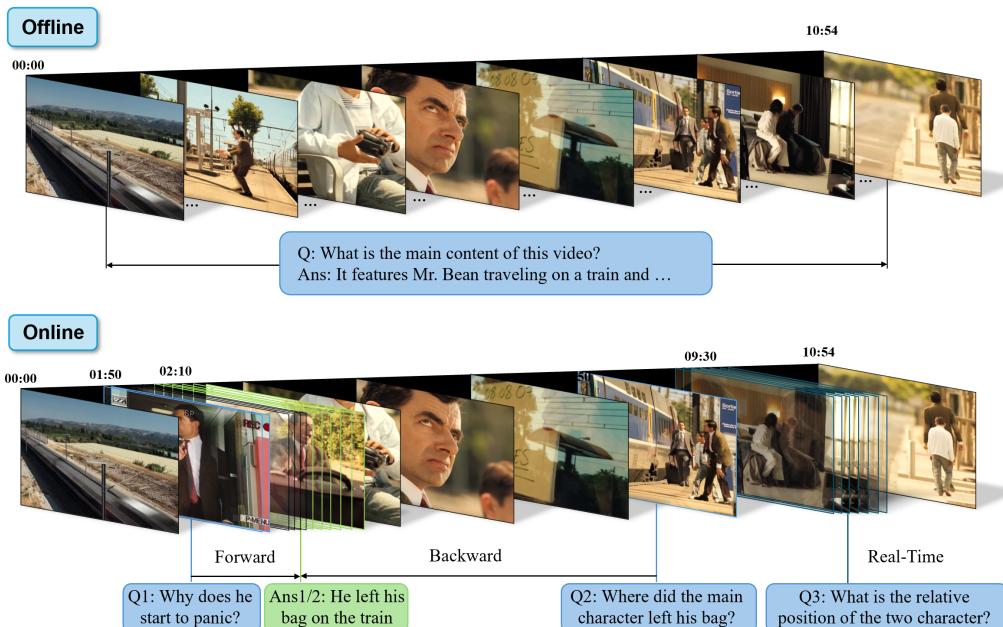


Figure 1. A demonstrative comparison between offline and online video understanding [6]. Offline video understanding focuses on answering questions based on the entirety of a video. In contrast, online video understanding involves posing queries about the context of a video at intermediate points, demanding the ability to trace back past information, perceive ongoing events, and adapt to continuous input.

Abstract

001 *Integrating past information and adapting to continuous*
002 *video input are pivotal for human-level video understand-*
003 *ing. Current benchmarks, however, focus on coarse-*
004 *grained, video-level question-answering in offline settings,*
005 *limiting real-time processing and adaptability for practical*
006 *applications. To this end, we introduce **OVBench** (Online-*
007 *Video-Benchmark), which assesses online video under-*
008 *standing through three modes: (1) **Backward Tracing**,*
009 *(2) **Real-Time Visual Perception**, and (3) **Forward Ac-***
010 *tive Responding*. OVBench consists of 12 tasks, compris-

011 *ing about 2,800 meta-annotations with fine-grained, event-*
012 *level timestamps paired with 858 videos across 10 do-*
013 *mains, encompassing egocentric activities, virtual gaming*
014 *worlds, and cinematic scenes. To minimize bias, we employ*
015 *automated generation pipelines and human annotation for*
016 *meticulous curation. We design an effective problem gener-*
017 *ation and evaluation pipeline based on these high-quality*
018 *samples and densely query Video-LLMs across the video*
019 *streaming timeline. Extensive evaluations of nine Video-*
020 *LLMs reveal that despite rapid advancements and improv-*
021 *ing performance on traditional benchmarks, existing mod-*
022 *els struggle with online video understanding. Our compre-*

023 *hensive evaluation reveals that the best-performing mod-*
024 *els still have a significant gap compared to human agents*
025 *in online video understanding. We anticipate that OVBench*
026 *will guide the development of Video-LLMs towards practi-*
027 *cal real-world applications and inspire future research in*
028 *online video understanding. Our benchmark and code can*
029 *be accessed at <https://github.com/JoeLeelyf/OVBench>.*

030 1. Introduction

031 Large Vision Language Models (LVLMs) [29, 36, 48, 60]
032 and Video-LLMs [25, 35, 58] have shown remarkable
033 progress, achieving impressive scores on existing bench-
034 marks [12, 13, 26]. Recent works like VideoLLM-online [6]
035 and Flash-VStream [59] have taken initial steps toward
036 J.A.R.V.I.S.-like real-world video assistants by combining
037 pre-trained vision encoders [41] and LLMs [10, 47]. How-
038 ever, a critical question remains: *How far are current state-*
039 *of-the-art models from achieving human-level online video*
040 *understanding?*

041 Despite the existence of dozens of evaluation bench-
042 marks in video understanding, there remains a signifi-
043 cant domain gap between these evaluations and real-world
044 video understanding tasks. Early evaluations [19, 54, 56]
045 primarily derive from video understanding and retrieval
046 datasets [3, 55], and assess models’ capability through
047 coarse-grained QAs, such as “*Q: Who is dancing? A:* Man”.
048 These QAs predominantly focus on short videos
049 with fixed question types and lack temporal indispensabil-
050 ity [12]. Subsequent works [13, 26, 62] attempt to ad-
051 dress these limitations by extending video temporal length
052 and incorporating more diverse tasks and video sources.
053 E.T.Bench [32] advances this further by exploring inher-
054 ent temporal information in videos and evaluating fine-
055 grained temporal event detection capabilities. However, all
056 works mentioned above are constrained to offline settings,
057 where models can access all video frames when respond-
058 ing to queries. While models that perform well on these
059 benchmarks demonstrate impressive capabilities in coarse-
060 grained video understanding, a significant gap between their
061 capability scope and the requirements for a real-world assis-
062 tant or an independent agent still exists.

063 Recent benchmarks VStream-QA [59] and Streaming-
064 Bench [28] introduce streaming understanding to the com-
065 munity. While VStream-QA offers an initial exploration
066 with limited video sources selected from Ego4d [23] and
067 MovieNet [17], StreamingBench provides an evaluation of
068 Video-LLMs’ in streaming scenarios at a larger scale. How-
069 ever, StreamingBench’s three major categories primarily fo-
070 cus on past visual inputs for real-time responses, offering an
071 incomplete view of streaming perception. We propose that
072 effective online video understanding requires simultaneous
073 capabilities to **trace back past information, perceive the**

074 **going-on, and adapt to the continuous input** simultane-
075 ously. Unlike StreamingBench’s emphasis on streaming
076 query format, our approach evaluates Video-LLMs’ ability
077 to find temporal visual clues from ongoing input, allowing
078 models to wait for sufficient evidence before responding.
079 We term this approach the **Video Chain-of-Time** thinking
080 process (Figure 3), analogous to Chain-of-Thought reason-
081 ing in LLMs [50].

082 We introduce **OVBench (Online-Video-Benchmark)** to
083 evaluate Video-LLMs’ online video understanding capa-
084 bilities. The benchmark comprises 858 videos from di-
085 verse sources, including curated datasets and web videos,
086 spanning 10 major domains (Sports, Video Games, 3D
087 scanning, etc.) with durations ranging from minutes to
088 half an hour. Using a hybrid approach combining semi-
089 automated MLLM generation and human curation, we cre-
090 ated 2814 high-quality samples (**Meta-Annotations**) with
091 precise event timestamps. These Meta-Annotations are or-
092 ganized into 12 tasks across three categories: **Backward**
093 **Tracing, Real-Time Visual Perception, and Forward Active**
094 **Responding**, reflecting the human video understanding
095 process illustrated in Fig. 3.

096 Building on the human-reviewed meta-annotations, we
097 develop an evaluation pipeline that queries Video-LLMs
098 densely along temporal axes to simulate continuous infor-
099 mation processing. For **Backward Tracing** and **Real-Time**
100 **Visual Perception**, we adopt multiple-choice evaluation,
101 converting videos into segments from start to query time
102 to accommodate offline models. With this approach, we
103 explore the potential of explicitly leveraging state-of-the-
104 art models in offline settings for online video understand-
105 ing. We evaluated 9 Video-LLMs, including proprietary
106 models GPT-4o [36] and Gemini-1.5-Pro [46], alongside
107 six recent open-source MLLMs like QWen2-VL [48] and
108 LLaVA-OneVision [24]. Despite their strong offline per-
109 formance, these models struggle with online-style queries
110 (e.g., *What is happening now?*), showing a significant gap
111 from human performance.

112 Further experiments with streaming models like Flash-
113 VStream [59] reveal even larger performance gaps,
114 with Flash-VStream achieving only 47.2% accuracy on
115 Real-Time Visual Perception compared to LLaVA-NeXT-
116 Video’s [61] 64.7% with equivalent parameters.

117 Our proposed **Forward Active Responding** setting in-
118 troduces the first evaluation framework specifically de-
119 signed for online video understanding, requiring models
120 to adapt responses to ongoing visual input continuously.
121 We evaluated the offline models using a novel multiple-
122 triggering pipeline that enables comparative assessment of
123 offline models’ potential in online video understanding and
124 implements score-based metrics for comprehensive evalua-
125 tion.

<p>126 2. Related Works</p> <p>127 Video Large Language Models. Video Large Language 128 Models (VLLMs) can process a video by treating it as a 129 sequence of video frames. Projects like VideoChat [25], 130 Video-LLaMA [58], and Video-ChatGPT [35] project 131 the CLIP-ViT [42] embeddings of selected video frames 132 through a Multi-Layer Perceptron (MLP) projector into the 133 LLM embedding space, then concatenate these embeddings 134 with text embeddings for enhanced video understanding. 135 However, the context length of MLLMs limits their effec- 136 tiveness in understanding long videos [25, 35], as longer 137 videos require more frames and a longer context length. To 138 address this limitation, two major approaches have been de- 139 veloped: compressing video features and selecting critical 140 frames.</p> <p>141 In the realm of feature compression, Chat-UniVi [22] 142 merges similar visual tokens through clustering techniques. 143 MovieChat [44] and MA-LLM [16] employ a memory 144 bank to store a fixed number of video tokens by iter- 145 atively merging the most similar tokens. ST-LLM [31] and 146 MovieChat [44] reduce video tokens to 32 using a pre- 147 trained Q-Former from BLIP2 [11]. LLaMA-VID [27] 148 takes a more radical approach, compressing each frame into 149 a content token and a context token.</p> <p>150 On the other hand, frame selection methods aim to 151 identify the most representative frames. VideoStreaming 152 [40] utilizes a small LLM to select critical video clips, 153 while FlashVstream [59] employs a clustering method to 154 choose representative frames for high-resolution process- 155 ing. LongVU [43] leverages question embeddings to se- 156 lect question-related frames, thereby enhancing video un- 157 derstanding.</p> <p>158 Benchmarks for Video Understanding. Traditional video 159 benchmarks, e.g., MSVD-QA [54], MSRVTT-QA [54], 160 and ActivityNet-QA [56], predominantly consist of short 161 videos, typically ranging from 1 to 2 minutes in duration. 162 These datasets are meticulously annotated with correspond- 163 ing questions and ground truth answers. GPT-4 [36] is em- 164 ployed to assess the accuracy of the answers by comparing 165 them against the provided questions and ground truth re- 166 sponses. However, these benchmarks primarily focus on 167 evaluating short, static video scenes. Hence, new bench- 168 marks designed to test causal and temporal understanding, 169 e.g., NExT-QA [53], TemporalBench [4], and AutoEval- 170 Video [8] are proposed.</p> <p>171 To gauge the capabilities of models on long-duration 172 videos, benchmarks like EgoSchema [34] covering over 173 5,000 egocentric videos with an average length of 3 min- 174 utes have been introduced. In contrast, Video-MME [13], 175 LVbench [49], and LongVideoBench [52] feature videos 176 spanning from 20 minutes to over an hour, evaluating 177 a broad spectrum of video understanding capabilities. 178 HourVideo [5] stands out with egocentric videos extending</p>	<p>179 180 181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197</p> <p>198 199 200 201 202 203 204</p> <p>205 206 207 208 209 210 211 212 213 214 215 216</p> <p>217 218 219 220 221 222 223 224 225</p>
<p>up to 2 hours, accompanied by more than 12,976 multiple-choice questions. Unlike these offline video benchmarks, our proposed OVBench is designed to evaluate online, interactive video understanding.</p> <p>Online Video Understanding. Traditional offline video understanding methodologies primarily focus on accessing entire video sequences to facilitate prediction tasks. Conversely, online video understanding demands models to process video streams sequentially, making decisions based on current and past information. This approach is particularly well-suited for scenarios where future data is unavailable, such as in embodied intelligence, autonomous driving, and augmented reality applications. Among online video understanding methods, FlashVStream [59] employs a clustering method to select representative frames, enabling MLLMs for real-time interactions. LIVE [6] introduces a comprehensive framework for learning in video streams, which includes a training objective, data generation schema, and an inference pipeline tailored for online video understanding.</p> <p>3. OVBench</p> <p>In this section, we present the construction process of our OVBench. We start with a detailed introduction to the three different modes of online video understanding, followed by a comprehensive description of the data collection and annotation procedures. A statistical report of our proposed benchmark is displayed at the end of this section.</p> <p>3.1. Online Video Understanding Mode Taxonomy</p> <p>Online video understanding aims to equip real-world, always-on agents with the ability to receive and process video inputs continuously, which closely mimics the human visual perception process. We categorize online video understanding into three distinct problem-solving modes: (1) Backward Tracing, (2) Real-Time Visual Perception, and (3) Forward Active Responding. Given a user-provided text query Q_{t_0} at the current time t_0 and a streaming video input $X_{(-\infty, +\infty)}$, these modes are formally defined as follows:</p> <ul style="list-style-type: none"> 1. Backward Tracing: $R_{t_0} = P(Q_{t_0}, X_{(-\infty, -T]})$ <p>2. Real-Time Visual Perception:</p> $R_{t_0} = P(Q_{t_0}, X_{(-T, t_0]})$ <p>3. Forward Active Responding:</p> $R_{(t_0, +\infty]} = P(Q_{t_0}, X_{(t_0, +\infty)})$ <p>in which T represents a threshold that defines the boundary for recent times, and R denotes the model's response. The first two modes, <i>Backward Tracing</i> and <i>Real-Time Visual Perception</i>, involve collecting visual information from</p>	

226 past and current timeframes respectively, and are expected
 227 to give immediate responses. In contrast, *Forward Active*
 228 *Responding* requires the model to withhold a response until
 229 sufficient future information becomes available to ensure
 230 a confident answer. Based on these distinctions, we have
 231 meticulously designed tasks tailored to each mode to ef-
 232 fectively evaluate the performance of Video-LLMs across
 233 these diverse capabilities.

234 3.1.1. Backward Tracing

235 Memory, particularly long-term memory, is a crucial aspect
 236 of human intelligence. In video understanding systems,
 237 this capability involves recalling and reasoning about past
 238 events. We focus on the following three tasks to evaluate
 239 this capability:

- 240 1. **[EPM] Episodic Memory:** Backtrack and retrieve key
 241 moments from past video inputs.
- 242 2. **[ASI] Action Sequence Identification:** Identify the cor-
 243 rect sequence of human actions in the video streams.
- 244 3. **[HLD] Hallucination Detection:** Ask questions irrele-
 245 vant to existing video inputs.

246 3.1.2. Real-Time Visual Perception

247 Accurate real-time perception of visual content is crucial, as
 248 actions undertaken in the present shape future outcomes. In
 249 various real-world scenarios, immediate and precise under-
 250 standing of ongoing visual inputs is essential. We propose
 251 six critical categories that constitute the foundational capa-
 252 bilities for effective real-time visual perception:

- 253 1. **[STU] Spatial Understanding.** Reason over the spa-
 254 tial relationships between objects occurring in nearby
 255 frames.
- 256 2. **[OJR] Object Recognition.** Recognize the objects ap-
 257 pearing in the current frames.
- 258 3. **[ATR] Attribute Recognition.** Identify the characteris-
 259 tics or properties of objects, such as color, texture, and
 260 size.
- 261 4. **[ACR] Action Recognition.** Recognize and interpret
 262 the actions being performed by individuals in the current
 263 frame.
- 264 5. **[OCR] Optical Character Recognition.** Recognize
 265 and interpret characters that appear within the frame.
- 266 6. **[FTP] Future Prediction.** Forecast the most probable
 267 subsequent phase of the current scene, including changes
 268 in object states, actions, and other dynamic elements.

269 3.1.3. Forward Active Responding

270 Transitioning from passive reception to active perception is
 271 essential for advanced video understanding systems. Ex-
 272 isting benchmarks primarily focus on the aforementioned
 273 two understanding modes, where Video-LLMs are required
 274 to respond immediately based on available information.
 275 In contrast, we introduce the *Forward Active Responding*
 276 mode, which allows the model to adjust its responses based



277 Figure 2. Examples of each task in OV-Bench. The 14 tasks are
 278 categorized into three different kinds of perceiving modes in
 279 online video understanding: **Backward Tracing**, **Real-Time Visual**
 280 **Perception**, and **Forward Active Responding**.

281 on forthcoming visual inputs. We devise four task dimensions to evaluate the models' active responding abilities:

- 282 1. **[REC] Repetition Event Count.** Respond when a repetitive event occurs again, including both high-frequency repetitive actions over short durations and semantically long-term occurrences of certain events.
- 283 2. **[SSR] Sequential Steps Recognition.** Respond when a certain procedure or sequence of actions has transitioned to another stage.
- 284 3. **[CRR] Clues Reveal Responding.** Delay responding until sufficient information or clues are provided.

285 3.2. Benchmark Construction

286 Under the taxonomy guidelines above, we make our first step by collecting video data and annotations from existing datasets and crawling data from the web to increase diver-

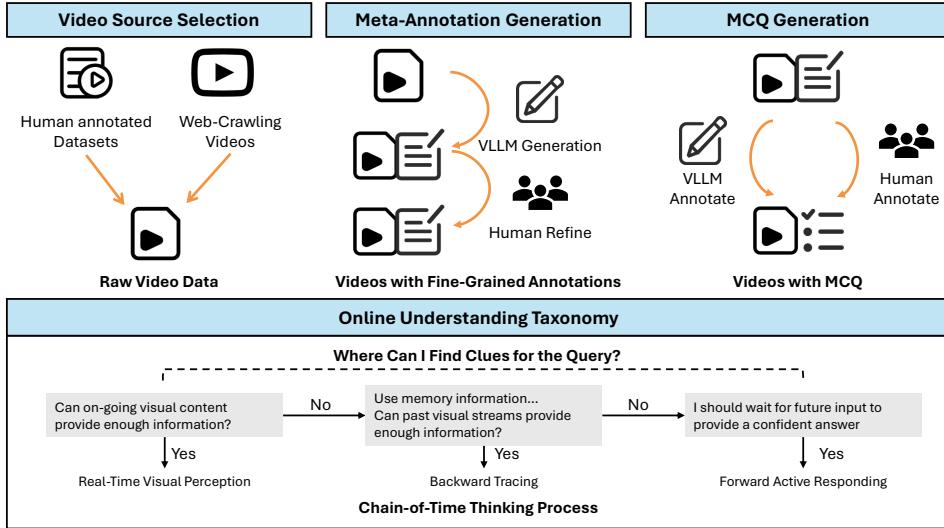


Figure 3. **Generation pipeline of OV-Bench.** Within public annotations, data is carefully filtered and relevant multiple-choice QAs are auto-generated. The effective system prompt and efficient answer prompt are employed to guide MLLMs toward precise outputs. The Video-LLMs we use to annotate videos are GPT-4o and Gemini-1.5 Pro.

sity. As our proposed evaluation pipeline highly relies on the accurate timestamp annotations of the referred events in the constructed prompt, the scarcity of event-level timestamps in existing datasets [51][33][39] promotes the design of our highly efficient meta-data generation pipeline 3. Raw annotations with coarse timestamps are then refined by humans to ensure accuracy. Our final questions and options for evaluation are constructed using our rule-based pipeline based on these human-refined meta-annotations. All QA samples undergo manual inspection before being included in the final test set.

3.2.1. Video and Annotation Collection

Video Source Selection. We follow existing benchmarks [32][26] by exploiting high-quality customized video datasets, and enrich our diversity by utilizing self-crawling videos from different domains. (1) **Human-annotated**

Video Dataset. Our main consideration for utilizing organized datasets is to alleviate the labor-intensive source video collection process. Specifically, we include QA-Ego4D[2] and OpenEQA[1] for the [EPM] task, STAR[51], YouCook2[63], CrossTask[65], HiREST[57], and COIN[45] for the [ASI] task, Perception-Test[39] and Thumos[21][14] for the [REC] task, COIN[45] for the [SSR] task, MovieNet[17] for the [CRR] task, and Ego4D[23] for tasks under *Real-Time Visual Perception*. All samples are selected from val or test sets to avoid potential data leakage. (2) **Web-crawling Videos.** To further extend the diversity of our benchmark, we follow the existing practice [12][28] of crawling source videos from YouTube.

Meta-Annotations Collection. We employ three approaches to collect our meta-annotations which contain event-level timestamps: (1) **Existing Annotation Re-**

purposing. For human-annotated datasets with accurate event-level timestamps [2][45][23], we explicitly take advantage of these labels and reconstruct them to our final prompt. (3) **Semi-Automatic Generation.** For datasets that provide video-level QA pairs without complete temporal localization, including [33][51][39][21][14], we prompt temporal-sensitive Video-LLMs like Gemini-1.5[46] to provide coarse-grain timestamps which fit the event referred in question and answer. For tasks under the *Real-Time Visual Perception* scenario, timestamps are given during our automatic QA construction process, which will be illustrated in 3.2.2. We then perform meticulously inspect all collected source videos and the corresponding meta-annotations to ensure precision.

3.2.2. Prompt Generation

Question and Answer Generation. Besides carefully selecting QA pairs from existing datasets to fit into our proposed tasks, we also adopt a highly efficient automatic question and answer generation pipeline, particularly for the *Real-Time Visual Perception* scenario. We randomly sample short clips from original long-form videos and then leverage GPT-4o[18] to select potential candidates and construct questions and corresponding answers using human-refined prompts. Human-proposed questions are also adopted as a part of these tasks to alleviate possible LLM preferences. For the novel [CRR] task, even the strongest Video-LLMs/MLLMs like Geimini-1.5-Pro struggle to construct desired problems. Volunteers are then recruited to provide QA pairs under our guidance.

Options Generation and Selection. We adopt multiple-choice questions as testing forms for *Backward Tracing* and *Real-Time Visual Perception* scenarios. However, as re-

326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358



Figure 4. **Left** Queries Temporal Distribution in OVBench. **Center** Linguistic Characteristics of Text Queries. **Right** Video category distribution of OVBench.

vealed in [7], the naively designed options of a multi-choice form query can cause information leakage about answers. We propose to generate options using a carefully designed rule-based and visually grounded transformation of correct answers, bringing misleading information from original videos to increase difficulty. Specifically, we prompt Video-LLMs with original QA pairs and corresponding video clips to generate visual-related options. A careful human review is then conducted to further ensure the options' effectiveness. All options are shuffled after human review to avoid potential preference bias.

Prompting Offline-Models for Simulated Online Understanding. With the significant performance gap between main-streaming powerful offline Video-LLMs [46][37][48] and existing online models [6][59], one natural question is made: *Is it effective to prompt offline models directly for online video understanding?* For the *Real-Time Visual Perception* setting, we make human curation to the original question to include implies about the real-time query scenarios, for example, by using sentence patterns like *What is/What am I* or containing words like *Now/Currently*. We made another intuitive attempt to prompt offline models to solve tasks under our novel *Forward Active Responding* scenario, which asks for a continuous adapting capability. Specifically, we devise a multiple-triggering densely query and evaluation pipeline, allowing the model to decide whether existing information has provided enough clues for answering the user's query.

3.3. Datasets Statistics

OVBench consists of 858 videos spanning 10 major domains, including Sports, Video Games, and 3D Scanning, among others. The video durations range from a few minutes to half an hour, with the average query timepoint being 428.89 seconds. Figure 4 **Left** illustrates the duration distribution of the queries within OVBench. The benchmark includes 2,814 question-answer (QA) pairs, featuring a large number of multiple-choice questions and a smaller set of open-ended questions. The number of options for the multiple-choice questions varies between 2 and 5, rather than being fixed at four. The distribution of video category

is visualized in Figure 4 **Right**.

4. Experiments

This section presents comprehensive experiments and in-depth analyses of OVBench.

4.1. Models and Evaluation Strategies

We evaluate four existing types of models: (1) Offline Multimodal Models, including GPT-4o [38], Gemini-1.5-Pro [46], Qwen2-VL [48], LLaVA-NeXT-Video [30], LLaVA-OneVision [29], InternVL-V2 [9] and LongVU [43], (2) Online Multimodal Models, including Flash-VStream-7B [59] and Videollm-Online[6] (3) Blind LLMs, including GPT-4-turbo [36]. (4) Human Agents. To ensure a fair comparison of model performance, we adhere to the principle of consistency by maintaining the same number of frames or frames per second (fps) across all models.

Considering the limitations on input video length for existing offline Video-LLMs, we adopt specialized video input methods tailored to such models. Specifically, we segment the video into clips based on the timestamps of the questions. For instance, for a question Q_i posed at timestamp t_i , we extract the video clip $\text{Video}[0 : t_i]$ as the visual input. This approach simulates a streaming question-answering scenario in online video understanding.

4.2. Main Results

Table 1 reports the performance of eleven models under different settings on OVBench, including the *Real-Time Visual Perception*, *Backward Tracing*, and *Forward Active Responding*. Our evaluation brings several important findings, as follows:

Offline Video-LLMs' video understanding capabilities can be effectively transferred to real-time video understanding. The results demonstrate that offline Video-LLMs, despite being designed for offline processing, perform competitively in *Real-Time Visual Perception* tasks. This suggests that the advanced video comprehension abilities developed in offline settings are transferable and can enhance performance in certain online scenarios, thereby

Model	# Frames	Real-Time Visual Perception						Backward Tracing				Forward Active Responding				Overall Avg.		
		OCR	ACR	ATR	STU	FTP	OJR	Avg.	EPM	ASI	HLD	Avg.	REC	SSR	CRR	Avg.	Overall Avg.	
Human																		
Human Agents																		
Blind LLMs																		
GPT-4-turbo																		
Proprietary Multimodal Models-Offline																		
Gemini 1.5 Pro	1fps	87.23	65.12	78.48	57.18	68.32	66.66	70.83	57.58	71.62	58.34	62.51	37.48	70.37	73.61	60.49	66.67	
GPT-4o	64	73.15	68.81	75.78	56.18	77.23	62.51	68.94	55.22	75.68	40.65	57.15	29.45	77.21	71.53	59.40	63.05	
Open-source Multimodal Models-Offline																		
Qwen2-VL-72B	64	77.18	61.47	76.72	51.68	76.24	64.62	67.98	52.19	75.68	68.02	65.30	36.41	75.50	63.19	58.37	63.88	
Qwen2-VL-7B	64	70.74	52.29	66.38	43.26	71.29	63.23	61.20	47.14	65.54	38.57	50.42	30.92	60.97	51.39	47.76	53.13	
InternVL-V2-8B	64	73.15	58.72	66.15	50.00	70.30	60.53	63.81	48.48	60.14	25.88	44.83	31.32	60.11	53.47	48.30	52.31	
LLaVA-NeXT-Video-7B	64	74.50	60.55	68.11	51.68	72.28	64.68	65.63	55.89	65.54	6.09	42.51	30.52	69.80	45.14	48.49	52.21	
LLaVA-OneVision-7B	64	70.74	59.63	73.28	47.19	71.29	63.64	64.63	56.23	56.08	16.74	43.02	27.58	63.82	38.19	43.20	50.28	
LongVU-7B	1fps	55.70	54.13	62.07	42.98	68.32	60.14	57.89	50.17	65.31	3.55	39.68	19.14	60.11	63.89	47.71	48.43	
Open-source Multimodal Models-Online																		
Flash-VStream-7B	1fps	40.94	51.38	44.90	38.52	62.38	44.97	47.18	42.76	38.51	3.05	28.11	15.8	40.15	65.97	40.64	38.64	
VideoLLM-online-8B	2fps	8.05	23.85	12.07	14.04	45.54	21.20	20.79	22.22	18.80	12.18	17.73	-	-	-	-	-	

Table 1. **Detailed evaluation results on OV-Bench.** † refers to using “low” resolution. To enhance the challenge of the questions by increasing the time interval between the question and the clues, the question time for [EPM] and [ASI] in the table is uniformly placed at the end of the video. [HLD] are evaluated in the standard manner following the streaming mode. For **Forward Active Responding**, accuracy-based evaluation metrics are utilized in this table. More high-quality [ASI] tasks are supplemented, further enhancing differentiation.

437 partially bridging the gap between offline and online video
438 understanding.

439 **Current Video-LLMs lack temporal prioritization**
440 when handling VQA tasks. Existing Video-LLMs do
441 not prioritize real-time temporal information when answer-
442 ing questions, leading to an inability to accurately locate
443 the correct scene when multiple misleading scenes match-
444 ing the question appear in the video stream, as shown in
445 Fig 2. Even the best current proprietary models achieve only
446 57.18% and 68.81% on [STU] and [ACR] tasks, respec-
447 tively, which represents a significant gap compared to Hu-
448 man Agents.

449 **A powerful LLM backbone is the key to achieving**
450 **high-performance video understanding.** As shown in Ta-
451 ble 1, the Qwen2-VL-72B model significantly outperforms
452 its smaller counterpart, Qwen2-VL-7B, across all evalua-
453 ted metrics. The larger model’s superior architecture al-
454 lows it to excel in *Real-Time Visual Perception* with an av-
455 erage score of 67.98%, compared to 61.20% for the 7B
456 variant. Additionally, Qwen2-VL-72B demonstrates en-
457 hanced capability in *Backward Tracing* tasks, achieving a
458 score of 68.02%, notably higher than the 54.15% scored
459 by Qwen2-VL-7B. These results highlight the importance
460 of a robust and powerful LLM backbone in effectively pro-
461 cessing and understanding complex video data, underscor-
462 ing the need for more powerful architectures to achieve high
463 performance in video understanding tasks.

464 **Hallucinations are prevalent in Video-LLMs.** The
465 [HLD] in Table 1 measures hallucinations in Video-
466 LLMs [20], indicating that hallucinations are a significant
467 issue, particularly in open-source and online models. Propri-
468 etary models like Gemini 1.5 Pro perform better in man-

aging hallucinations, yet there remains a notable gap com-
469 pared to human performance(58.34% vs. 91.37%). This
470 problem arises due to the models’ inability to fully com-
471 prehend complex visual and temporal contexts, leading to
472 errors in interpretation and response. Addressing halluci-
473 nations is crucial for improving the reliability and accuracy of
474 Video-LLMs in real-world applications.

4.3. Comparison between online Video-LLMs and 476 offline Video-LLMs

477 Models like **Gemini 1.5 Pro** and **Qwen2-VL-72B**, repre-
478 sentative of offline Video-LLMs, demonstrate strong per-
479 formance across various tasks, as shown in Fig 7. Specifically,
480 Gemini 1.5 Pro achieves the highest average score among
481 these models. This superior performance suggests that of-
482 fline models, despite not being designed for online or real-
483 time processing, can effectively comprehend and process
484 complex visual information when provided with sufficient
485 computational resources and pre-processing time. Their ar-
486 chitectures typically allow for processing the entire video
487 sequence holistically, leveraging global context and detailed
488 temporal information, which enhances their temporal un-
489 derstanding and reasoning capabilities.

490 In contrast, **Flash-VStream-7B**, representing online
491 Video-LLMs, shows comparatively lower performance in
492 real-time perception tasks compared to offline models. This
493 model is designed to process video in a streaming man-
494 nner, handling inputs frame by frame with strict latency con-
495 straints to achieve real-time responsiveness. The perfor-
496 mance gap highlights a potential trade-off between real-
497 time processing capabilities and the depth of visual under-
498 standing.

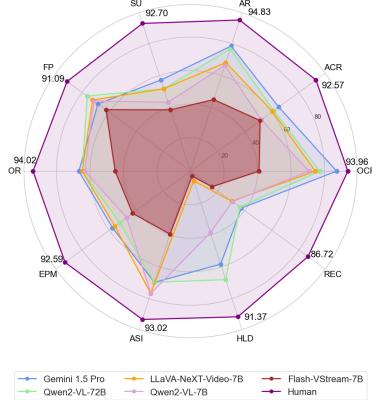


Figure 5. **Performance comparison between online Video-LLMs and offline Video-LLMs.** The figure illustrates the average scores of different models on the OVbench in real-time visual perception tasks.

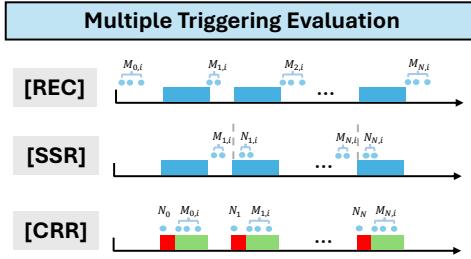


Figure 6. **Multiple triggering evaluation pipeline of prompt offline models for online video understanding.** Offline Video-LLMs are densely queried along the temporal axes to make independent decisions of whether existing visual content provide enough clues for answering.

4.4. Forward Active Responding

We include our evaluation pipeline design for our proposed *Forward Active Responding*. While our high-quality human-annotated queries and clues lay an ideal testbed for future real-world online understanding models, existing naively designed online video models usually collapse in our evaluation process. We made our initial attempts to leverage our multiple-triggering query pipeline to prompt offline VideoLLMs to perform online video understanding thinking schema and further explore their potential in always-on visual perception.

Evaluation Pipeline and Metrics. As illustrated in Fig.6, We propose to query the Video-LLMs densely along the temporal axes, particularly around the interested events'. Our main concerns are twofold: 1) Encourage models' timely finding of the right clues, and 2) Avoid any possible hallucination before the right clue appears. For the

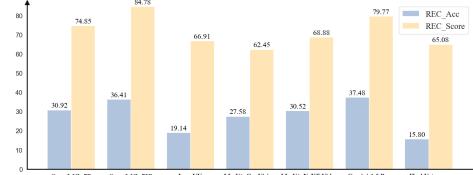


Figure 7. **Performance of models in Repetition Event Count and Clue Reveal Responding task.** Evaluation results are given using the average accuracy and our proposed metrics.

[REC] task, larger counting numbers are awarded. Based on this, we proposed our designed scoring metrics for the three tasks in the *Forward Active Responding*. Given N the total times of events, and i/j the corresponding trigger times, our evaluation metrics is formed as below:

1. Repetition Event Count:

$$Score = \sum_i^N e^{n/5} \cdot \sum_{j=0}^3 acc_j^{\frac{1}{i+1}}$$

2. Sequential Steps Recognition:

$$Score = Acc \cdot \sum_{i=0}^2 p_i \cdot 2^{1-i}$$

3. Clues Reveal Responding:

$$Score = Acc \cdot \sum_{i=0}^3 p_i \cdot 2^{1-i}$$

Offline Models for Online Video Understanding. Despite their promising performance on the *Backward-Tracing* and *Real-Time Visual Perception*, in which the models are given full information for making confident responses, our preliminary results show that even state-of-the-art offline models like Geimini-1.5-Pro, fails to capture the linguistic information of ongoing querying, showing limited understanding of online video content.

5. Conclusion and Future Work

In this work, we introduced OVbench, a comprehensive benchmark designed to assess online video understanding capabilities of Video-LLMs across three critical modes: *Backward Tracing*, *Real-Time Visual Perception*, and *Forward Active Responding*. We anticipate that OVbench will serve as a valuable resource for the research community, guiding the development of Video-LLMs toward practical, real-world applications. By highlighting current limitations and providing a platform for rigorous evaluation, we hope to inspire future research dedicated to advancing online video understanding and achieving human-level comprehension in artificial intelligence systems.

- 549 **References**
- 550 [1] Xiaohan Zhang, Pranav Putta Sriram Yenamandra Mikael
551 Henaff, Sneha Silwal Paul Mcvay Oleksandr MakSYMets Sergio Arnaud Karmesh Yadav Qiyang Li Ben Newman Mohit
552 Sharma Vincent Berges Shiqi Zhang Pulkit Agrawal Yonatan Bisk Dhruv Batra Minal Kalakrishnan Franziska Meier
553 Chris Paxton Sasha Sax Aravind Rajeswaran Arjun Majumdar, Anurag Ajay. OpenEQA: Embodied Question Answering
554 in the Era of Foundation Models. In *CVPR*, 2024. 5, 4
- 555 [2] Leonard Bärmann and Alex Waibel. Where did i leave my
556 keys?-episodic-memory-based question answering on ego-
557 centric videos. In *Proceedings of the IEEE/CVF Conference
558 on Computer Vision and Pattern Recognition*, pages 1560–
559 1568, 2022. 5, 4
- 560 [3] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem,
561 and Juan Carlos Niebles. Activitynet: A large-scale video
562 benchmark for human activity understanding. In *Proceedings
563 of the ieee conference on computer vision and pattern
564 recognition*, pages 961–970, 2015. 2
- 565 [4] Mu Cai, Reuben Tan, Jianrui Zhang, Bocheng Zou, Kai
566 Zhang, Feng Yao, Fangrui Zhu, Jing Gu, Yiwu Zhong,
567 Yuzhang Shang, Yao Dou, Jaden Park, Jianfeng Gao,
568 Yong Jae Lee, and Jianwei Yang. Temporalbench: Bench-
569 marking fine-grained temporal understanding for multimodal
570 video models, 2024. 3
- 571 [5] Keshigeyan Chandrasegaran, Agrim Gupta, Lea M Hadzic,
572 Taran Kota, Jimming He, Cristóbal Eyzaguirre, Zane Du-
573 rante, Manling Li, Jiajun Wu, and Li Fei-Fei. Hourvideo:
574 1-hour video-language understanding. *arXiv preprint arXiv:2411.04998*, 2024. 3
- 575 [6] Joya Chen, Zhaoyang Lv, Shiwei Wu, Kevin Qinghong Lin,
576 Chenan Song, Difei Gao, Jia-Wei Liu, Ziteng Gao, Dongxing
577 Mao, and Mike Zheng Shou. Videollm-online: Online video
578 large language model for streaming video. In *CVPR*, 2024.
579 1, 2, 3, 6
- 580 [7] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang
581 Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao,
582 Dahua Lin, et al. Are we on the right way for evaluating large
583 vision-language models? *arXiv preprint arXiv:2403.20330*,
584 2024. 6
- 585 [8] Xiuyuan Chen, Yuan Lin, Yuchen Zhang, and Weiran Huang.
586 Autoeval-video: An automatic benchmark for assessing
587 large vision language models in open-ended video question
588 answering, 2024. 3
- 589 [9] Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen,
590 Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu,
591 Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng
592 Dai. Internvl: Scaling up vision foundation models and
593 aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*, 2023. 6
- 594 [10] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao
595 Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao
596 Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source
597 chatbot impressing gpt-4 with 90%* chatgpt quality. See
598 <https://vicuna.lmsys.org> (accessed 14 April 2023), 2023.
599 2
- 600 [11] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat
601 Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale
602 Fung, and Steven Hoi. Instructblip: Towards general-
603 purpose vision-language models with instruction tuning.
604 *arXiv preprint arXiv:2305.06500*, 2023. 3
- 605 [12] Xinyu Fang, Kangrui Mao, Haodong Duan, Xiangyu Zhao,
606 Yining Li, Dahua Lin, and Kai Chen. Mmbench-video: A
607 long-form multi-shot benchmark for holistic video under-
608 standing. *arXiv preprint arXiv:2406.14515*, 2024. 2, 5
- 609 [13] Chaoyou Fu, Yuhan Dai, Yondong Luo, Lei Li, Shuhuai Ren,
610 Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen,
611 Mengdan Zhang, et al. Video-mme: The first-ever compre-
612 hensive evaluation benchmark of multi-modal llms in video
613 analysis. *arXiv preprint arXiv:2405.21075*, 2024. 2, 3
- 614 [14] Alex Gorban, Haroon Idrees, Yu-Gang Jiang, A Roshan Zamir,
615 Ivan Laptev, Mubarak Shah, and Rahul Sukthankar. Thumos
616 challenge: Action recognition with a large number
617 of classes, 2015. 5, 4
- 618 [15] Kristen Grauman, Andrew Westbury, Eugene Byrne,
619 Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson
620 Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d:
621 Around the world in 3,000 hours of egocentric video. In *Pro-
622 ceedings of the IEEE/CVF Conference on Computer Vision
623 and Pattern Recognition*, pages 18995–19012, 2022. 4
- 624 [16] Bo He, Hengduo Li, Young Kyun Jang, Menglin Jia, Xuefei
625 Cao, Ashish Shah, Abhinav Shrivastava, and Ser-Nam Lim.
626 Ma-lmm: Memory-augmented large multimodal model for
627 long-term video understanding supplementary material. 3
- 628 [17] Qingqiu Huang, Yu Xiong, Anyi Rao, Jiaze Wang, and
629 Dahua Lin. Movienet: A holistic dataset for movie under-
630 standing. In *Computer Vision–ECCV 2020: 16th European
631 Conference, Glasgow, UK, August 23–28, 2020, Proceedings,
632 Part IV 16*, pages 709–727. Springer, 2020. 2, 5, 4
- 633 [18] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perel-
634 man, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda,
635 Alan Hayes, Alec Radford, et al. Gpt-4o system card.
636 *arXiv preprint arXiv:2410.21276*, 2024. 5
- 637 [19] Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and
638 Gunhee Kim. Tgif-qa: Toward spatio-temporal reasoning in
639 visual question answering. In *Proceedings of the IEEE con-
640 ference on computer vision and pattern recognition*, pages
641 2758–2766, 2017. 2
- 642 [20] Shaoliang Chen Jingjing Chen Yu-Gang Jiang Ji-
643 acheng Zhang, Yang Jiao. Eventhallusion: Diagnosing
644 event hallucinations in video llms. *arXiv preprint arXiv:
645 2409.16597*, 2024. 7
- 646 [21] Yu-Gang Jiang, Jingen Liu, A Roshan Zamir, George
647 Toderici, Ivan Laptev, Mubarak Shah, and Rahul Sukthankar.
648 Thumos challenge: Action recognition with a large number
649 of classes, 2014. 5, 4
- 650 [22] Peng Jin, Ryuichi Takanobu, Caiwan Zhang, Xiaochun Cao,
651 and Li Yuan. Chat-univi: Unified visual representation em-
652 powers large language models with image and video under-
653 standing. *arXiv preprint arXiv:2311.08046*, 2023. 3
- 654 [23] Eugene Byrne Zachary Chavis Antonino Furnari Rohit Gird-
655 har Jackson Hamburger Hao Jiang Miao Liu Xingyu Liu
656 et al. Kristen Grauman, Andrew Westbury. Ego4d: Around
657

- the world in 3,000 hours of egocentric video. In *CVPR*, pages 18995–19012, 2022. 2, 5
- [24] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 2
- [25] KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023. 2, 3
- [26] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22195–22206, 2024. 2, 5
- [27] Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language models. 2024. 3
- [28] Junming Lin, Zheng Fang, Chi Chen, Zihao Wan, Fuwen Luo, Peng Li, Yang Liu, and Maosong Sun. Streamingbench: Assessing the gap for mllms to achieve streaming video understanding. *arXiv preprint arXiv:2411.03628*, 2024. 2, 5
- [29] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. 2, 6
- [30] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. 6
- [31] Ruyang Liu, Chen Li, Haoran Tang, Yixiao Ge, Ying Shan, and Ge Li. St-lm: Large language models are effective temporal learners. In *European Conference on Computer Vision*, pages 1–18. Springer, 2025. 3
- [32] Ye Liu, Zongyang Ma, Zhongang Qi, Yang Wu, Chang Wen Chen, and Ying Shan. E.t. bench: Towards open-ended event-level video-language understanding. In *Neural Information Processing Systems (NeurIPS)*, 2024. 2, 5, 4
- [33] Arjun Majumdar, Anurag Ajay, Xiaohan Zhang, Pranav Putta, Sriram Yenamandra, Mikael Henaff, Sneha Silwal, Paul Mcvay, Oleksandr Maksymets, Sergio Arnaud, et al. Openeqa: Embodied question answering in the era of foundation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16488–16498, 2024. 5
- [34] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding, 2023. 3
- [35] Salman Khan Muhammad Maaz, Hanoona Rasheed and Fahad Khan. Video-chatpt: Towards detailed video understanding via large vision and language models. *ArXiv 2306.05424*, 2023. 2, 3
- [36] OpenAI. Gpt-4 technical report, 2023. Technical report. 2, 3, 6
- [37] OpenAI. Hello gpt-4o. <https://openai.com/index/hello-gpt-4o>, 2024. 6
- [38] OpenAI. Hello gpt-4o, 2024. <https://openai.com/index/hello-gpt-4o>. 6
- [39] Viorica Patraucean, Lucas Smaira, Ankush Gupta, Adria Re-casens, Larisa Markeeva, Dylan Banarse, Skanda Koppula, Mateusz Malinowski, Yi Yang, Carl Doersch, et al. Perception test: A diagnostic benchmark for multimodal video models. *Advances in Neural Information Processing Systems*, 36, 2024. 5, 4
- [40] Rui Qian, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Shuangrui Ding, Dahua Lin, and Jiaqi Wang. Streaming long video understanding with large language models. *arXiv preprint arXiv:2405.16009*, 2024. 3
- [41] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 2
- [42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 3
- [43] Xiaoqian Shen, Yunyang Xiong, Changsheng Zhao, Lemeng Wu, Jun Chen, Chenchen Zhu, Zechun Liu, Fanyi Xiao, Balakrishnan Varadarajan, Florian Bordes, Zhuang Liu, Hu Xu, Hyunwoo J. Kim, Bilge Soran, Raghuraman Krishnamoorti, Mohamed Elhoseiny, and Vikas Chandra. Longvu: Spatiotemporal adaptive compression for long video-language understanding. *arXiv preprint arXiv:2410.17434*, 2024. 3, 6
- [44] Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Xun Guo, Tian Ye, Yan Lu, Jenq-Neng Hwang, and Gaoang Wang. Moviechat: From dense token to sparse memory for long video understanding, 2023. 3
- [45] Yansong Tang, Dajun Ding, Yongming Rao, Yu Zheng, Danyang Zhang, Lili Zhao, Jiwen Lu, and Jie Zhou. Coin: A large-scale dataset for comprehensive instructional video analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1207–1216, 2019. 5, 4
- [46] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024. 2, 5, 6
- [47] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 2
- [48] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 2, 6
- [49] Weihan Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Xiaotao Gu, Shiyu Huang, Bin Xu, Yuxiao Dong, et al. Lvbench: An extreme long video understanding benchmark. *arXiv preprint arXiv:2406.08035*, 2024. 3
- [50] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny

- 777 Zhou. Chain-of-thought prompting elicits reasoning in large
778 language models, 2023. 2
- 779 [51] Bo Wu, Shoubin Yu, Zhenfang Chen, Joshua B Tenenbaum,
780 and Chuang Gan. Star: A benchmark for situated reasoning
781 in real-world videos. *arXiv preprint arXiv:2405.09711*,
782 2024. 5, 4
- 783 [52] Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li.
784 Longvideobench: A benchmark for long-context interleaved
785 video-language understanding, 2024. 3
- 786 [53] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua.
787 Next-qa: Next phase of question-answering to explaining
788 temporal actions. In *Proceedings of the IEEE/CVF conference*
789 on computer vision and pattern recognition, pages
790 9777–9786, 2021. 3
- 791 [54] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang,
792 Xiangnan He, and Yueteng Zhuang. Video question answering
793 via gradually refined attention over appearance and motion.
794 In *Proceedings of the 25th ACM international conference on*
795 *Multimedia*, pages 1645–1653, 2017. 2, 3
- 796 [55] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large
797 video description dataset for bridging video and language. In
798 *Proceedings of the IEEE conference on computer vision and*
799 *pattern recognition*, pages 5288–5296, 2016. 2
- 800 [56] Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yuet-
801 ing Zhuang, and Dacheng Tao. Activitynet-qa: A dataset for
802 understanding complex web videos via question answering.
803 In *Proceedings of the AAAI Conference on Artificial Intelli-*
804 *gence*, pages 9127–9134, 2019. 2, 3
- 805 [57] Abhay Zala, Jaemin Cho, Satwik Kottur, Xilun Chen, Bar-
806 las Oguz, Yashar Mehdad, and Mohit Bansal. Hierarchical
807 video-moment retrieval and step-captioning. In *Proceedings*
808 of the IEEE/CVF Conference on Computer Vision and Pat-
809 tern Recognition, pages 23056–23065, 2023. 5, 4
- 810 [58] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An
811 instruction-tuned audio-visual language model for video un-
812 derstanding. *arXiv preprint arXiv:2306.02858*, 2023. 2, 3
- 813 [59] Haoji Zhang, Yiqin Wang, Yansong Tang, Yong Liu, Jiashi
814 Feng, Jifeng Dai, and Xiaojie Jin. Flash-vstream: Memory-
815 based real-time understanding for long video streams, 2024.
816 2, 3, 6
- 817 [60] Pan Zhang, Xiaoyi Dong, Bin Wang, Yuhang Cao, Chao Xu,
818 Linke Ouyang, Zhiyuan Zhao, Shuangrui Ding, Songyang
819 Zhang, Haodong Duan, Wenwei Zhang, Hang Yan, Xinyue
820 Zhang, Wei Li, Jingwen Li, Kai Chen, Conghui He,
821 Xingcheng Zhang, Yu Qiao, Dahua Lin, and Jiaqi Wang.
822 Internlm-xcomposer: A vision-language large model for ad-
823 vanced text-image comprehension and composition, 2023. 2
- 824 [61] Yuanhan Zhang, Bo Li, haotian Liu, Yong jae Lee, Liangke
825 Gui, Di Fu, Jiashi Feng, Ziwei Liu, and Chunyuan Li. Llava-
826 next: A strong zero-shot video understanding model, 2024.
827 2
- 828 [62] Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao, Xi
829 Yang, Yongping Xiong, Bo Zhang, Tiejun Huang, and Zheng
830 Liu. Mlvu: A comprehensive benchmark for multi-task
831 long video understanding. *arXiv preprint arXiv:2406.04264*,
832 2024. 2
- 833 [63] Luowei Zhou, Chenliang Xu, and Jason Corso. Towards
834 automatic learning of procedures from web instructional
835 videos. In *Proceedings of the AAAI Conference on Artificial
836 Intelligence*, 2018. 5
- 837 [64] Luowei Zhou, Chenliang Xu, and Jason Corso. Towards
838 automatic learning of procedures from web instructional
839 videos. In *Proceedings of the AAAI Conference on Artificial
840 Intelligence*, 2018. 4
- 841 [65] Dimitri Zhukov, Jean-Baptiste Alayrac, Ramazan Gokberk
842 Cinbis, David Fouhey, Ivan Laptev, and Josef Sivic. Cross-
843 task weakly supervised learning from instructional videos.
844 In *Proceedings of the IEEE/CVF Conference on Computer
845 Vision and Pattern Recognition*, pages 3537–3545, 2019. 5,
846 4