

StreamReady: Learning *What* to Answer and *When* in Long Streaming Videos

Anonymous CVPR submission

Paper ID ****

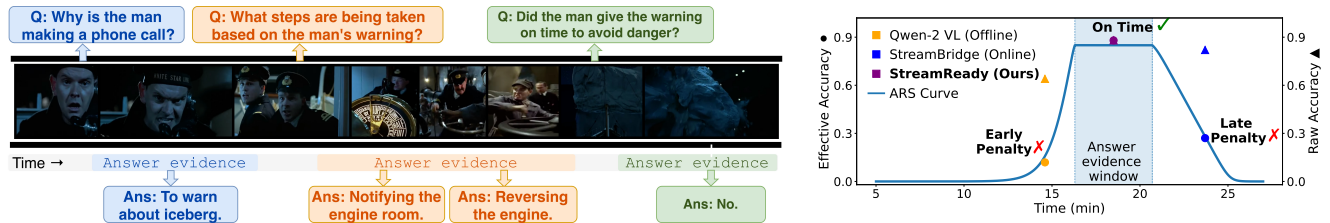


Figure 1. **Readiness-aware streaming video understanding.** *Left:* In proactive streaming settings, questions can precede their supporting evidence, requiring the model to monitor the evolving video and answer once the relevant cues appear. *Right:* Under our readiness-aware formulation, effective accuracy jointly reflects answer correctness and timing via the Answer Readiness Score (ARS). Although all models achieve similar raw accuracy on this example, ARS reveals sharp performance drops for early (hallucinatory) or late (delayed) answers. In contrast, StreamReady responds within the evidence window, preserving high effective accuracy by answering at the appropriate moment.

Abstract

001 Streaming video understanding often involves time-sensitive
 002 scenarios where models need to answer exactly when the
 003 supporting visual evidence appears: answering before the
 004 evidence reflects speculation, answering after it has passed
 005 reduces real-time utility. To capture this behavior, we in-
 006 troduce a readiness-aware formulation of streaming video
 007 understanding with the **Answer Readiness Score (ARS)**,
 008 a timing-aware objective with asymmetric early and late
 009 penalties. When combined with correctness, ARS defines an
 010 effective accuracy that measures not just whether a model
 011 is right, but whether it answers at the appropriate moment.
 012 Building on this formulation, we introduce **StreamReady**,
 013 a framework to unify temporal reasoning with on-time an-
 014 swering through a lightweight readiness mechanism that
 015 decides if sufficient evidence has been observed before re-
 016 sponding. To evaluate this capability, we further introduce
 017 **ProReady-QA**, a benchmark with annotated answer evi-
 018 dence windows and proactive multi-turn questions across
 019 local and global contexts. StreamReady achieves super-
 020 ior performance on ProReady-QA, and consistently out-
 021 performs prior methods across eight additional streaming
 022 and offline long-video benchmarks, demonstrating robust
 023 and broadly generalizable video understanding capability.

024 1. Introduction

026 Multimodal Large Language Models (MLLMs) have sig-
 027 nificantly advanced video understanding across diverse do-
 028 mains [1, 2, 5, 20, 23, 38], particularly on short clips.

029 However, their performance drops on long videos due to
 030 difficulties of reasoning over extended temporal context
 031 [59]. While recent efforts extend MLLMs to longer videos
 032 [3, 8, 14, 36, 39, 57], these models still operate offline with
 033 full-video access during inference. In contrast, streaming
 034 video understanding represents an online variant of long-
 035 video reasoning, where frames arrive sequentially, evidence
 036 may appear before or after a question, and the model must
 037 operate without ever seeing the full video. This online ca-
 038 pability is essential for real-world settings such as surveil-
 039 lance, sports analytics, robotics, and assistive systems that
 040 demand timely and context-aware reasoning.

041 Building on advances in long-video reasoning, recent
 042 studies have begun exploring streaming video understand-
 043 ing [9, 18, 32, 37, 43, 47, 52]. However, most existing
 044 works focus on past-dependent (causal) reasoning where
 045 evidence is already available at question time, making an-
 046 swer correctness the main objective; leaving timing largely
 047 unexplored. In contrast, Many real-world scenarios require
 048 future-dependent (proactive) reasoning, where the question
 049 appears *before* the supporting evidence, requiring the model
 050 to watch the unfolding video to determine when enough
 051 information has appeared (Figure 1, *left*). In such proac-
 052 tive settings, models must prioritize answer timing as much
 053 as correctness: responding too early, even if correct, indi-
 054 cates unsupported speculation; responding too late causes
 055 unwanted delay. Developing this ability to identify the right
 056 moment to respond based on supporting evidence is there-
 057 fore essential for truly effective streaming understanding.

058 Recently models have begun exploring such proactive

059	behavior by deferring responses through auxiliary MLLMs	112
060	[37] or prompt-based cues [51], though at the cost of non-	113
061	determinism or added compute. To evaluate this behav-	114
062	ior, benchmarks [26, 27, 41, 42] include proactive scenar-	
063	ios to encourage models to wait until relevant information	
064	appears. However, they lack annotated answer evidence du-	
065	rations, making it difficult to verify if responses are given at	
066	an appropriate time. Consequently, these developments of-	
067	fer only a partial view of timing behavior since models may	
068	wait, but they lack any criteria for determining whether their	
069	chosen answer time is supported by the actual evidence.	
070	To address these limitations, we formalize readiness-	
071	aware streaming video understanding, where the goal is not	
072	only to produce the correct answer, but to do so precisely	
073	<i>when</i> sufficient evidence appears. At its core is the Ans-	
074	wer Readiness Score (ARS) , a timing-aware evaluation	
075	metric with asymmetric penalties: a harsher early penalty	
076	discourages unsupported guesses before evidence, and a	
077	milder late penalty tolerates slight delays after the evidence	
078	ends (Figure 1 <i>right</i>). Together, these penalties yield an ef-	
079	fective accuracy that captures both answer correctness <i>and</i>	
080	timing. Building on this formulation, we propose Stream-	
081	Ready , a framework that unifies temporal reasoning with	
082	explicit answer-timing. Instead of relying on heavy auxil-	
083	iary models or prompt heuristics, StreamReady introduces	
084	a lightweight learnable readiness token within its reason-	
085	ing module, allowing the model to assess from its internal	
086	memory when sufficient evidence has appeared. A small	
087	readiness head monitors this token, prompting the model to	
088	answer only when appropriate, ensuring responses are both	
089	accurate and timely.	
090	To evaluate readiness-aware understanding, we intro-	
091	duce ProReady-QA , a benchmark designed for proactive	
092	scenarios with annotated answer evidence durations and	
093	tasks covering both local and global temporal contexts,	
094	enabling systematic evaluation under ARS. Our evalua-	
095	tions show that StreamReady effectively bridges the gap	
096	between raw and effective accuracy, outperforming exist-	
097	ing methods in readiness-aware streaming settings. Beyond	
098	ProReady-QA, StreamReady also achieves superior perfor-	
099	mance on other streaming benchmarks across proactive and	
100	non-proactive tasks, and generalizes well to offline long-	
101	video benchmarks. Together, our proposed formulation,	
102	method, and benchmark provide a unified foundation for	
103	advancing answer timing in streaming video understanding.	
104	Our main contributions are as follows:	
105	• We formalize readiness-aware streaming understanding	
106	and introduce the Answer Readiness Score (ARS) to	
107	jointly evaluate answer correctness and timing through	
108	asymmetric early and late penalties.	
109	• We propose StreamReady , a readiness-aware framework	
110	to integrate temporal reasoning with a readiness mecha-	
111	nism to decide evidence sufficiency before responding.	
	• We develop ProReady-QA , a benchmark with annotated	112
	answer evidence windows and proactive multi-turn ques-	113
	tions for evaluating timing behavior of streaming models.	114
	2. Related Works	115
	Offline Long Video Understanding with MLLMs. Long-	116
	video understanding aims to model extended temporal con-	117
	text in videos spanning from few minutes to hours. Prior	118
	works use memory-based [4, 8, 14, 36] or agent-based	119
	approaches [19, 30, 39, 57], with some adopting query-	120
	conditioned storage [3, 44]. However, these offline methods	121
	rebuild memory for every query and rely on full-video ac-	122
	cess, making them unsuitable for streaming settings where	123
	frames and questions arrive sequentially. Our work draws	124
	inspiration from query-aware conditioning but adapts it for	125
	streaming understanding without requiring memory reset.	126
	Streaming Video Understanding with MLLMs. Inspired	127
	by long-video frameworks, memory-based [31, 32, 37, 43,	128
	45, 50, 53] and retrieval-based approaches [9, 18, 47] ex-	129
	tend MLLMs to streaming by processing frames online and	130
	reusing past context. While effective for answer content,	131
	they lack mechanisms for deciding answer timing; particu-	132
	larly important in proactive scenarios where questions pre-	133
	cede the answer. Our readiness-aware design complements	134
	these approaches by adding explicit answer timing, avoid-	135
	ing premature speculation and unnecessary delays by re-	136
	sponding precisely when the evidence appears.	137
	Streaming Video Benchmarks. Existing streaming bench-	138
	marks [17, 24, 45, 50, 53] primarily support past-dependent	139
	question-answering, where timing has limited impact. More	140
	recent benchmarks include proactive scenarios [26, 27, 41,	141
	42], but remain limited to short clips and local context. In	142
	contrast, ProReady-QA supports proactive reasoning over	143
	long, continuous streams with both local and global multi-	144
	turn dependencies, enabling comprehensive evaluation of	145
	on-time answering in readiness-aware streaming.	146
	3. Method	147
	A fundamental requirement in streaming video understand-	148
	ing is not only <i>what</i> a model answers but also <i>when</i> it	149
	chooses to answer. We formalize this as readiness-aware	150
	streaming video understanding , where a model must pro-	151
	duce the correct answer at the appropriate moment, sup-	152
	ported by visual evidence. Unlike existing streaming setups	153
	that evaluate only correctness and therefore cannot distin-	154
	guish on-time answers from mistimed ones, our formulation	155
	explicitly accounts for both answer content <i>and</i> timing.	156
	3.1. Framework Overview	157
	Building on this formulation, we introduce StreamReady ,	158
	a readiness-aware framework that learns to determine the	159
	right moment to answer by monitoring the evolving video	160
	for supporting evidence. As the video unfolds, Stream-	161

162 Ready stores them in a hierarchical memory (§3.2),
 163 retrieves and reasons over temporally relevant context when
 164 a question appears (§3.3) and uses a lightweight readiness
 165 mechanism (§3.4) to decide whether sufficient evidence is
 166 present to answer. If ready, it answers immediately; oth-
 167 erwise, it continues observing until the required evidence
 168 appears. Figure 2 provides an overview of the framework.

169 3.2. Memory Storage

170 Streaming video understanding can benefit from both vi-
 171 sual and semantic history, where visual cues guide percep-
 172 tion and prior linguistic interactions provide context, en-
 173 abling efficient reuse of past information. Since streaming
 174 videos contain fine-grained details and substantial tempo-
 175 ral redundancy, an effective system must preserve key evi-
 176 dence while compacting redundant content. StreamReady
 177 achieves this through two complementary memories: a Vi-
 178 sual Memory Tree (\mathcal{M}_V) for multi-granular visual context,
 179 and a Contextual Memory Bank (\mathcal{M}_C) for long-range se-
 180 mantic dependencies across question–answer rounds.

181 **Visual Memory Tree (\mathcal{M}_V).** To efficiently represent long
 182 streaming videos, we maintain a multi-level Visual Mem-
 183 ory Tree that progressively abstracts incoming frames. The
 184 lowest level \mathcal{M}_{V1} stores the most recent frame embed-
 185 dings in a FIFO buffer, preserving short-term details. Once full,
 186 its raw frames are compressed into a compact centroid set
 187 $\mathcal{M}_{V2} = \{c_1, c_2, \dots, c_J\}$ via K-means clustering, forming a
 188 stable mid-level summary of the recent segment. As stream-
 189 ing continues, evicted frames f_o of \mathcal{M}_{V1} , update \mathcal{M}_{V2}
 190 through EMA-based clustering with decay factor α :

$$191 \quad c_j \leftarrow \begin{cases} (1 - \alpha)c_j + \alpha f_o, & \text{if } \text{sim}(f_o, c_j) \geq \tau_t, \\ \text{new centroid}, & \text{otherwise} \end{cases} \quad (1)$$

192 where, threshold τ_t tightens in stable scenes (favoring
 193 merges) and relaxes when novelty rises (allowing new clus-
 194 ters), keeping \mathcal{M}_{V2} compact yet adaptive. When \mathcal{M}_{V2}
 195 reaches capacity J or shows distributional drift (e.g., fre-
 196 quent new-cluster creation), its centroids are abstracted into
 197 a coarse prototype set $\mathcal{M}_{V3} = \{s_1, s_2, \dots, s_U\}$, also up-
 198 dated through EMA-based clustering to support fast contin-
 199 uous abstraction during stable scenes.

$$200 \quad s_u \leftarrow (1 - \alpha)s_u + \alpha \left(\frac{1}{|\mathcal{I}_u|} \sum_{j \in \mathcal{I}_u} c_j \right), \quad (2)$$

201 If centroids become heterogeneous (e.g., low mutual simi-
 202 larity or persistent novelty), a lightweight mini-K-means is
 203 triggered to realign prototypes, improving coherence with-
 204 out full re-clustering. This hierarchical design yields a com-
 205 pact multi-granular memory that preserves long-range fine-
 206 grained details while reducing redundancy, enabling effi-
 207 cient retrieval for query-aware reasoning.

208 **Contextual Memory Bank (\mathcal{M}_C).** Beyond visual evi-
 209 dence, many streaming questions depend on earlier linguis-

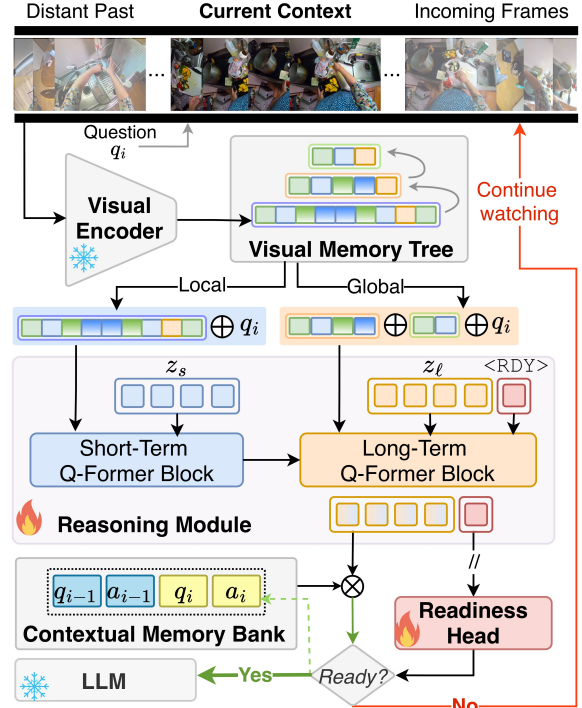


Figure 2. **Framework Overview.** StreamReady encodes stream-
 ing videos into a visual memory tree and reasons through short
 and long-term branches. A learnable $\langle \text{RDY} \rangle$ token, guided by a readi-
 ness head, gates the reasoning output until sufficient evidence is
 observed. Once ready, the long-term representation, enriched with
 contextual information from past QA pairs, is sent to the LLM for
 answering, enabling readiness-aware streaming behavior.

210 tic interactions. We support this with a Contextual Mem-
 211 ory Bank \mathcal{M}_C where each entry stores the question embed-
 212 ding (q_i) and the learned representation (a_i) that was used
 213 to generate its answer, forming a lightweight semantic his-
 214 tory, that provides a complementary view to \mathcal{M}_V , enabling
 215 efficient context reuse for multi-turn reasoning.

216 3.3. Query Aware Reasoning

217 Once a question q_i arrives, the model transitions from pas-
 218 sive encoding to active, query-aware retrieval and reason-
 219 ing. Because evidence in streaming video is distributed un-
 220 evenly over time, the model must reason over both recent
 221 and older fine-grained context and abstracted summaries.
 222 To support this, we use a dual-branch Q-Former [3], consist-
 223 ing of a short-term branch $\mathcal{Q}_s(\cdot)$ that operates on \mathcal{M}_{V1} ,
 224 and a long-term branch $\mathcal{Q}_l(\cdot)$, that operates on $\mathcal{M}_{V2}, \mathcal{M}_{V3}$.
 225 **Short-term Reasoning.** To capture short-range query-
 226 relevant cues, the short-term branch $\mathcal{Q}_s(\cdot)$, attends to the
 227 raw frames in \mathcal{M}_{V1} , together with q_i , producing the short-
 228 term learned representation z_s :

$$229 \quad z_s = \mathcal{Q}_s(\text{Concat}[\mathcal{M}_{V1}, q_i]), \quad (3)$$

230 **Long-term Reasoning.** In long streaming videos, relevant
 231 evidence is often buried in earlier moments, but only a small

portion might be relevant to any given query. To efficiently reason over useful long-range evidence, we perform coarse-to-fine query-aware retrieval over the higher levels of \mathcal{M}_V .

We first score all prototypes in \mathcal{M}_{V_3} and select the top-K high-level regions most likely to contain relevant evidence:

$$\mathcal{M}_{V_3}^t = \text{Top-K}((W_s q_i)^\top \mathcal{M}_{V_3}), \quad (4)$$

This acts as a semantic map lookup that highlights the most relevant video regions. The prototype scores are normalized (via softmax) to stabilize routing and sharpen focus on plausible evidence. For each selected prototype s , we then gather its associated centroids c from \mathcal{M}_{V_2} to form a candidate pool C_0 , refine their relevance score, and select the top- m fine-grained slots:

$$\mathcal{M}_{V_2}^t = \underset{c \in C_0}{\text{Top-m}} ((W_c q_i)^\top c); \quad C_0 = \bigcup_{s \in \mathcal{M}_{V_3}^{(s)}} \mathcal{M}_{V_2}^{(s)}, \quad (5)$$

Unlike the prototype selection stage, we avoid normalizing centroid scores since sharp ranking is essential for isolating specific evidence rather than broad regions. The retrieved prototypes $\mathcal{M}_{V_3}^t$, centroids $\mathcal{M}_{V_2}^t$, and the short-term learned representation z_s are then passed to the long-term branch $\mathcal{Q}_\ell(\cdot)$ for reasoning:

$$z_\ell = \mathcal{Q}_\ell(\text{Concat}[\mathcal{M}_{V_3}^t, \mathcal{M}_{V_2}^t, q_i], z_s), \quad (6)$$

This two-stage retrieval mirrors episodic recall: prototypes provide coarse temporal anchors, and centroids supply the fine-grained details. Combining them with z_s enables efficient long-term reasoning that is grounded in both broad and recent evidence, producing a focused representation.

Contextual Reasoning. To effectively reuse past semantic context, we perform a contextual reasoning step over the contextual memory \mathcal{M}_C . The current question embedding is matched with stored previous question embeddings \mathcal{M}_C to identify similar past QA interactions. A soft-gating mechanism selects the most relevant entries, and their answer representations are fused into the long-term visual feature z_ℓ through a lightweight cross-attention layer. This complements visual reasoning by incorporating prior semantic knowledge for coherent multi-turn understanding.

3.4. Readiness Mechanism

The reasoning modules determine *what* to answer, but, readiness-aware streaming understanding also requires deciding *when* to answer. To support this, we introduce a learnable <RDY> token and a Readiness Head.

Monitoring Readiness. We append the <RDY> token to the long-term reasoning representation z_ℓ within $\mathcal{Q}_\ell(\cdot)$ which already learns evidence from the visual memory tree, exposing the <RDY> token to the same evolving evidence used for answer reasoning. This ensures readiness decisions are grounded in the same representations that drive answer generation. Before the relevant evidence appears,

retrieval remains weakly aligned with the question, producing a diffused, low-confidence reasoning state. As supporting evidence arrives and retrieval becomes sharper and more question-consistent; the representation shifts towards an answer-bearing state. By residing within $\mathcal{Q}_\ell(\cdot)$, the <RDY> token naturally learns to track this transition and encode how prepared the model is to answer at a given time.

A lightweight Readiness Head, monitors this token and outputs a readiness score $R_{\text{pred}} \in [0, 1]$ at each timestep. At inference, the model triggers the LLM to respond only when this score exceeds a threshold; otherwise, it continues observing future frames. When the model is ready to answer, the fused representation of z_ℓ and \mathcal{M}_C is passed to the LLM for response generation, and this fused representation is stored in \mathcal{M}_C as the answer representation a_i for the current question q_i . This gating enforces readiness-aware behavior by preventing premature guesses and unnecessary delays, ensuring responses occur precisely when the model judges the evidence to be sufficient.

Learning the Readiness Signal. Since the model must function online during inference but training has full-video access, we leverage this to construct weak pseudo-supervision without ever requiring ground-truth evidence timestamps. We estimate likely evidence locations by measuring similarity between the learned representation z_ℓ and centroid-level memory \mathcal{M}_{V_2} across time. \mathcal{M}_{V_2} offers a suitable balance between detail and compactness because its centroids capture fine-grained cues without the redundancy of \mathcal{M}_{V_1} or heavy abstraction of \mathcal{M}_{V_3} . High-similarity centroids define a pseudo-positive temporal region P , while low-similarity ones define a pseudo-negative region N . The readiness mechanism is trained to assign higher readiness to P than to N through a pairwise contrastive loss:

$$\mathcal{L}_{ctr} = -\log \sigma (R_{pred}(t^+) - R_{pred}(t^-)), \quad (7)$$

where $t^+ \in P, t^- \in N$. Because \mathcal{L}_{ctr} alone can produce noisy, unstable readiness signals, we add a mild temporal coherence regularizer in the final objective:

$$\mathcal{L}_{rdy} = \mathcal{L}_{ctr} + \lambda_{reg} \|\nabla_t R_{pred}(t)\|_1, \quad (8)$$

where $\|\cdot\|_1$ denotes the L1 norm. Crucially, \mathcal{L}_{rdy} updates only the Readiness Head and the <RDY> token; gradients are stopped from the rest of the reasoning module so that the reasoning pathways learn *what* to answer through its standard video-text loss, while the readiness mechanism independently learns *when* to answer through the timing loss.

4. ProReady-QA Benchmark and Evaluation

To support readiness-aware streaming understanding evaluation, we introduce **ProReady-QA**, a benchmark designed specifically for long-duration streaming videos with proactive multi-turn questions and explicitly annotated answer evidence windows across local and global context (Table



Figure 3. Examples of each task in ProReady-QA. Here, the question and answer frames are color-coded.

330 1). ProReady-QA enables the first systematic evaluation of
 331 both answer correctness and temporal appropriateness, pro-
 332 viding a dedicated testbed for studying answer timing be-
 333 havior in streaming video models.

334 4.1. Task Definition

335 ProReady-QA spans five proactive reasoning tasks (Figure
 336 3) requiring models to track future evidence and make tem-
 337 porally aware decisions: (1) **Sequential Steps Recogni-**
 338 **tion (SSR)**: detect process transitions; (2) **Repetitive Event**
 339 **Count (REC)**: count recurring events; (3) **Clues Reveal**
 340 **Responding (CRR)**: answer once certain evidence appears;
 341 (4) **Causal Trigger Detection (CTD)**: detect cause-effect
 342 events; (5) **Goal-State Detection (GSD)**: identify goal
 343 completion. The first three build on prior proactive tasks
 344 [26, 27], while GSD and CTD extend proactive reasoning
 345 to higher-level temporal and causal understanding.

346 4.2. Dataset Construction

347 ProReady-QA contains 10 one-hour Ego-4D [13] and 22
 348 half-hour MovieNet [16] videos, with 5k proactive QA pairs
 349 with annotated answer evidence windows timestamps.

350 **Sourcing Videos.** We build upon long videos from
 351 VStream-QA [53], which offer diverse activities and suf-
 352 ficient temporal depth to support future-dependent reason-
 353 ing. However, VStream-QA’s questions are entirely past-
 354 dependent, making them ill-suited for proactive scenario.

355 **Generating QA and Evidence Timestamp.** We remove
 356 all past-dependent questions from VStream-QA and create

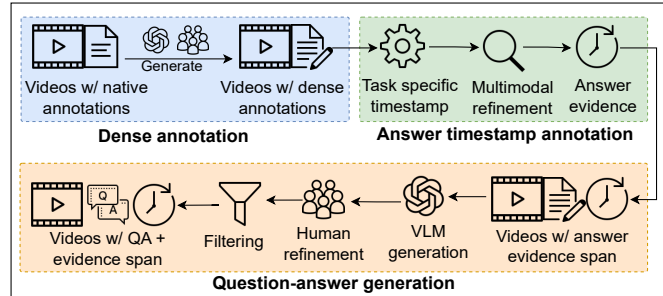


Figure 4. Generation pipeline of ProReady-QA.

Table 1. Comparison of ProReady-QA with prior benchmarks.

Benchmark	QA (K)	V Len. (min)	V Avg. (min)	Multi-turn QA		Ans. Evidence
				Local	Global	
Non-Proactive Streaming Reasoning Benchmarks						
ODVBench [50]	6.3	<1-3	1.2	×	×	×
StreamBench [45]	1.8	2-9	4.5	×	×	×
OVBench [17]	7.0	<1-3	1.5	×	×	×
VStream-QA [53]	3.5	30-60	40	×	×	×
Proactive Streaming Reasoning Benchmarks						
ProactiveVideoQA [41]	1.4	<1-2	2.1	×	×	×
OVOBench [26]	2.8	<1-30	7.1	×	×	×
Omni-MMI [42]	2.3	<1-12	5.4	✓	×	×
StreamingBench [27]	4.5	1-10	4.1	✓	×	×
ProReady-QA	5.0	30-60	40	✓	✓	✓

new multi-turn proactive questions whose answers rely on
 future frames, following a semi-automatic QA generation
 pipeline [26, 53] with human refinement. Complex tasks
 (CRR and CTD) include manually authored questions to
 capture nuanced temporal structure. Each QA pair is an-
 notated with precise evidence windows by aligning native
 annotations with multimodal cues (visual frames, subtitles,
 action boundaries) to identify first- and last-valid evidence
 timestamps. For tasks with extended events (SSR, REC),
 we define onset and end times to cover the full target ev-
 idence duration, yielding fine-grained temporal boundaries
 essential for timing evaluation. ProReady-QA also incorpo-
 rates local and global multi-turn dependencies, where later
 questions reference earlier entities or events, reflecting real-
 world long-horizon reasoning. Figure 4 illustrates the over-
 all QA generation process.

4.3. Evaluating Answer Readiness

To evaluate readiness-aware streaming performance, we in-
 troduce the **Answer Readiness Score (ARS)**, a timing-
 aware evaluation metric that penalizes answers given too
 early or too late relative to their evidence. When combined
 with accuracy (Acc), ARS produces an effective accuracy
 (Acc_e) that rewards predictions that are both correct and
 well-timed. For a set of N questions, ARS is defined as:

$$ARS = \frac{1}{N} \sum_{i=1}^N (EP_i \cdot LP_i); \quad Acc_e = Acc \times ARS \quad (9)$$

Table 2. **Performance comparison of readiness-aware streaming understanding on ProReady-QA.** †: Qwen-2-VL backbone. Metrics include accuracy (Acc.), Answer Readiness Score (ARS), effective accuracy (Acc_e). **Best** and **second-best** performances are highlighted.

Method	Size	SSR		CRR		REC		GSD		CTD		Average		
		Acc.	ARS	Acc.	ARS	Acc.	ARS	Acc.	ARS	Acc.	ARS	Acc.	ARS	Acc _e
<i>Offline Video MLLMs</i>														
InternVL2 [7]	8B	57.2	0.17	51.7	0.39	24.3	0.43	29.4	0.39	31.8	0.49	38.9	0.37	0.20
LLaVA-OneVision [21]	7B	71.8	0.48	58.2	0.53	22.3	0.45	39.4	0.24	34.8	0.18	45.3	0.38	0.29
Qwen-2-VL [38]	7B	67.9	0.32	53.1	0.39	20.7	0.31	35.1	0.52	30.3	0.28	41.4	0.34	0.20
LLaVA-NeXT-Video [55]	7B	67.4	0.47	58.3	0.41	24.3	0.42	40.3	0.31	21.8	0.38	42.4	0.40	0.31
MiniCPM-V 2.6 [49]	8B	62.7	0.56	54.2	0.62	24.5	0.35	32.7	0.43	23.4	0.35	39.5	0.46	0.23
HierarQ [3]	7B	67.4	0.32	55.8	0.45	28.7	0.44	45.2	0.38	32.8	0.43	46.0	0.40	0.27
<i>Online Video MLLMs</i>														
VideoLLM-online [6]	8B	53.2	0.30	51.1	0.38	18.3	0.12	30.5	0.34	24.9	0.44	29.6	0.32	0.18
Flash-VStream [53] †	7B	71.4	0.61	60.3	0.59	22.4	0.31	42.8	0.48	21.4	0.35	43.7	0.47	0.34
Dispider [32] †	7B	47.2	0.59	52.3	0.55	20.3	0.20	34.2	0.33	27.2	0.41	36.2	0.42	0.27
StreamForest [50] †	7B	65.8	0.41	57.4	0.51	24.4	0.24	49.2	0.29	32.3	0.40	45.8	0.37	0.27
StreamBridge [37] †	7B	<u>72.2</u>	<u>0.72</u>	59.7	<u>0.65</u>	31.9	0.43	60.3	<u>0.57</u>	41.4	<u>0.49</u>	<u>53.1</u>	<u>0.60</u>	<u>0.42</u>
ViSpeak [12] †	7B	67.5	0.61	55.2	0.54	25.3	0.27	48.2	0.49	27.3	0.23	44.7	0.43	0.31
InfiniPot-V [18] †	7B	69.4	0.50	<u>60.8</u>	0.52	<u>34.2</u>	<u>0.54</u>	58.3	0.42	37.2	0.35	52.0	0.47	0.36
StreamReady †	7B	74.3	0.78	63.3	0.73	39.6	0.68	61.2	0.68	43.5	0.59	56.4	0.69	0.53

where EP and LP represent early and late penalties. For each question, ProReady-QA provides a ground-truth evidence window $[t_s, t_e]$, where t_s marks when sufficient evidence first appears and t_e when it ceases to be valid. Given a model answers at time t_a , these intervals enable ARS to evaluate timing behavior through complementary early and late penalties aligned with our readiness-aware formulation. **Early Penalty (EP).** Answering *before* any supporting evidence appears is the most severe readiness failure, as it reflects speculation rather than observation. To discourage such behavior, the Early Penalty sharply decreases as the model answers earlier than the evidence onset t_s , scaled by the median evidence duration τ for consistency.

$$EP = \text{softmin} \left(1, 2 \sigma \left(\gamma_e \frac{t_a - t_s}{\tau + \epsilon} \right) \right) \quad (10)$$

Here, $\sigma(\cdot)$ is sigmoid, γ_e controls penalty sharpness, ϵ ensures numerical stability. $EP \rightarrow 1$ as t_a approaches t_s , while early answers yield lower scores.

Late Penalty (LP). After the evidence ends at t_e , delayed responses indicate hesitation rather than hallucination. To capture this milder form of readiness error, the Late Penalty gently decreases with delay, encouraging timely responses without overly penalizing slight delays:

$$LP = \text{softmin} \left(1, \text{softmax} \left(0, 1 - \gamma_\ell \frac{t_a - t_e}{\tau + \epsilon} \right) \right) \quad (11)$$

where γ_ℓ controls decay slope, keeping LP near 1 for minor delays and lower for prolonged ones. For questions

with multiple valid answers, we compute ARS separately for each turn and report their average.

5. Experiments

Datasets and Metrics. We evaluate StreamReady on ProReady-QA using accuracy, ARS and effective accuracy. To show generalization beyond readiness-aware streaming, we further evaluate on four streaming benchmarks: StreamingBench [27], OVOBench [26], OVBench [17], and VStream-QA [53]; and four offline long-video benchmarks: VideoMME [10], MLVU [58], MVBench [22], and EgoSchema [29], following each dataset’s official evaluation protocol and reporting accuracy.

Implementation Details. We use Qwen-2 VL [38] as the backbone model and initialize the dual-branch Q-Former using pretrained weights [3]. For fair timing evaluation, offline models are tested on ProReady-QA by truncating videos at the evidence end and recording when the correct answer first appears, while streaming models process frames sequentially with the prompt “Answer whenever you are ready” following [26, 27]. For offline long-video benchmarks, readiness and contextual reasoning are disabled. We set $\gamma_e = 6$, $\gamma_\ell = 1$ for ARS penalties and use a readiness threshold of 0.35 for LLM trigger following [12, 37, 40].

5.1. Results

Readiness-Aware Streaming Understanding. Table 2 shows that StreamReady achieves the highest accuracy and ARS across all five ProReady-QA tasks, surpassing the best model by $\sim 3\%$ in accuracy and $\sim 9\%$ in ARS on average,

Table 3. Performance comparison of streaming video understanding. †: Qwen-2-VL backbone. Accuracy is the reported metric.

(a) Benchmark: StreamingBench (Strm), and OVOBench (OVO).										(b) Benchmark: VStream-QA. Non-Proactive.			
Method	Size	Proactive: ✓		Proactive: ✗				Average		Method	Size	RE	RM
		Strm	OVO	Strm		OVO		Strm	OVO				
				Cont.	Fwd.	Real	Omni						
<i>Offline Video MLLMs</i>										<i>Online Video MLLMs</i>			
InternVL2 [7]	8B	32.4	45.4	63.7	35.8	60.7	44.0	44.0	50.7	7B	50.7	36.0	
LLaVA-OneVision [21]	7B	32.7	50.9	71.1	38.4	62.8	45.0	47.4	56.4	7B	56.4	49.4	
Qwen-2-VL [38]	7B	31.7	48.9	69.0	34.9	60.7	48.6	45.2	57.3	7B	57.3	53.1	
LongVU [34]	7B	-	48.5	-	-	57.4	39.5	-	55.8	7B	55.8	50.8	
LongVA [54]	7B	30.2	-	63.1	35.9	-	-	43.1	57.9	7B	57.9	51.4	
LLaVA-NeXT-Video [55]	7B	34.3	54.2	69.8	41.7	63.3	41.7	48.6	64.8	7B	64.8	57.2	
VITA 1.5 [11]	7B	27.4	53.5	52.3	33.1	63.5	41.5	37.6	-	-	-	-	
MiniCPM-V 2.6 [49]	8B	35.0	-	67.4	35.0	-	-	45.8	-	-	-	-	
HierarQ [3]	7B	35.1	48.3	69.7	44.4	67.3	48.3	49.7	54.6	-	-	-	
<i>Online Video MLLMs</i>										(c) Benchmark: OVBench. Non-Proactive.			
VideoLLM-online [6]	8B	26.6	-	36.0	28.5	20.8	17.7	30.4	41.9	7B	41.9	39.1	
Flash-VStream [53] †	7B	24.1	44.2	23.2	26.0	29.9	25.4	24.4	39.1	8B	39.1	30.4	
Dispider [32] †	7B	33.6	34.7	67.6	35.7	54.6	36.1	45.6	20.4	7B	20.4	48.7	
StreamForest [50] †	7B	-	53.5	77.3	-	61.2	52.0	-	48.7	8B	48.7	43.6	
ReKV [9] w/o offload	7B	30.7	69.1	37.4	-	-	-	-	49.5	7B	49.5	48.7	
StreamBridge [37] †	7B	32.6	48.4	77.0	24.1	71.3	68.1	44.6	48.7	7B	48.7	43.6	
ViSpeak [12] †	7B	43.9	54.3	70.4	61.6	66.3	57.5	58.6	59.4	8B	59.4	59.4	
InfiniPot-V [18] †	7B	-	47.9	76.4	-	65.9	47.6	-	53.8	7B	53.8	53.8	
TimeChat-Online [48]	7B	35.3	36.4	75.4	37.8	58.6	42.0	49.5	45.7	7B	45.7	45.7	
StreamAgent [46] †	7B	34.6	45.4	74.3	36.3	61.3	41.7	48.4	49.5	7B	49.5	49.5	
StreamReady †	7B	48.2	58.8	78.3	63.7	73.6	72.2	63.4	68.2	7B	63.9	63.9	

Table 4. Performance comparison of offline long-video understanding. †: Qwen-2-VL backbone. Accuracy is reported.

Method	Size	VidMME	MLVU	MVB	EgoSch
<i>Open-source Offline Video MLLMs</i>					
VideoChat-GPT [28]	7B	-	31.3	32.7	49.6
LLaMA-VID [25]	7B	33.2	41.9	38.5	38.5
InternVL2 [7]	8B	54.0	64.0	65.8	55.0
LongVA [54]	7B	52.6	56.3	51.3	46.7
LLaVA-OneVision [21]	7B	58.2	64.7	56.7	60.1
Qwen-2-VL [38]	7B	63.3	65.8	67.0	66.7
LongVU [34]	7B	60.6	65.4	66.9	67.6
LLaVA-Video [56]	7B	63.3	70.8	58.6	57.3
HierarQ [3]	7B	63.7	69.4	67.6	67.3
<i>Open-source Online Video MLLMs</i>					
MovieChat [35]	7B	38.2	25.8	55.1	53.5
Flash-VStream [53] †	7B	61.2	66.3	65.4	68.2
VideoChat-online [17]	4B	52.8	60.8	64.9	54.7
Dispider [32] †	7B	57.2	61.7	-	55.6
StreamForest [50] †	7B	61.4	70.0	70.2	-
ReKV [9] w/o offload	7B	-	68.5	-	60.7
InfiniPot-V [18] †	7B	62.8	65.8	-	65.6
VideoLLaMB [43]	7B	41.4	-	52.5	-
TimeChat-Online [48]	7B	62.5	65.4	-	-
StreamBridge [37] †	7B	64.4	69.6	64.4	66.9
ViSpeak [12] †	7B	55.0	54.1	54.1	-
StreamReady †	7B	65.8	71.3	71.8	70.4

435 with the largest ARS gains on REC, GSD, and CTD tasks.
 436 StreamReady’s readiness mechanism reduces mistimed re-
 437 sponses, producing tighter temporal alignment that directly
 438 lifts ARS. Additionally, its hierarchical visual memory ex-

439 poses the model to clearer evidence boundaries, making
 440 readiness estimation more reliable and leading to a smaller
 441 gap between raw and effective accuracy.

442 **Streaming video understanding.** Tables 3a–3c shows
 443 StreamReady consistently outperforms prior models,
 444 achieving up to ~5% gains on proactive tasks of stream-
 445 ing benchmarks. While the readiness mechanism improves
 446 proactive behavior by encouraging the model to wait for
 447 sufficient evidence, most of the accuracy improvements in
 448 these benchmarks come from StreamReady’s stronger evi-
 449 dence retrieval and long-horizon reasoning. Its visual mem-
 450 ory offers reliable access to relevant visual cues, and the
 451 contextual memory provides complementary semantic his-
 452 tory for multi-turn interactions. Combined with query-
 453 aware long-range reasoning, this design suppresses irrele-
 454 vant historical frames and emphasizes the segments most
 455 predictive of an answer, leading to more stable and accurate
 456 predictions even when timing itself is not evaluated.

457 **Offline long video understanding.** Table 4 shows that
 458 StreamReady also maintains strong performance in offline
 459 long-video understanding, outperforming prior models de-
 460 spite temporal awareness being irrelevant. Here, the benefit
 461 mainly comes from its ability to integrate long-range visual
 462 structure as the memory hierarchy offers compact yet ex-
 463 pressive summaries of extended video segments, allowing
 464 the model to reason accurately over long temporal spans.

Table 5. Ablation studies for each component.

Method	REC		GSD		CTD	
	Acc.	ARS	Acc.	ARS	Acc.	ARS
Contribution of each component						
Baseline	20.7	0.31	35.1	0.52	30.3	0.28
+ Trivial Reasoning [3]	28.7	0.44	48.2	0.53	34.2	0.36
+ Readiness Mechanism	28.4	0.50	47.9	0.58	34.3	0.42
+ Memory Storage	32.4	0.46	50.2	0.52	38.8	0.36
+ Query-aware Reasoning	39.4	0.48	60.9	0.53	43.6	0.39
+ Readiness Mechanism	39.6	0.68	61.2	0.68	43.5	0.59
Design choice of Readiness Mechanism						
Only MLP	39.5	0.52	60.9	0.55	42.9	0.42
LLM w/ Heuristic Triggers [51]	39.6	0.54	61.4	0.54	43.5	0.43
Auxilliary MLLM [37]	39.2	0.60	61.2	0.61	43.3	0.46
<RDY> + Head (Transformer)	39.5	0.69	61.6	0.68	43.1	0.58
<RDY> + Head (MLP)	39.6	0.68	61.2	0.68	43.5	0.59
Placement of <RDY> token						
Input of Q_s	39.1	0.31	60.9	0.51	43.1	0.18
Learned representation of Q_s	39.6	0.38	61.3	0.54	42.8	0.24
Input of Q_ℓ	39.4	0.54	61.1	0.62	43.6	0.49
Learned representation of Q_ℓ	39.6	0.68	61.2	0.68	43.5	0.59

465 5.2. Ablation

466 We ablate StreamReady’s readiness mechanism on three
467 challenging ProReady-QA tasks (REC, GSD, and CTD).
468 Additional architectural ablations are in Supplementary.

469 **Contribution of each component.** Table 5 (top) shows
470 that while adding basic reasoning [3] improves accuracy,
471 it provides little timing benefit. Incorporating the readiness
472 mechanism on it increases ARS modestly, but major gains
473 occur with our stronger memory and reasoning modules.
474 Together, these components enhance both accuracy and timing,
475 showing that improved evidence retrieval strengthens
476 not only *what* but also *when* the model answers.

477 **Design choice of readiness mechanism.** Table 5 (mid-
478 dle) shows that an MLP-based head is insufficient to judge
479 readiness, while LLM or MLLM-based readiness adds only
480 modest improvement due to weak coupling with reasoning.
481 Embedding the <RDY> token within the reasoning module
482 offers the best and most stable ARS improvements by di-
483 rectly accessing evolving evidence, with a lightweight MLP
484 head matching Transformer performance at lower cost.

485 **Design choice of <RDY> token placement.** Table 5 (bot-
486 tom) shows that placing <RDY> in the short-term branch
487 yields noisy timing signals due to its limited local context,
488 while positioning it as input to the long-term branch offers
489 moderate but unstable results as it only passively attends
490 to retrieved evidence. The best ARS is achieved by attach-
491 ing <RDY> to the learned long-term representation, where it
492 co-evolves with query-aligned evidence, enabling the most
493 reliable readiness detection.

494 5.3. Analysis

495 **Penalty Sharpness for Early and Late Responses.** Fig-
496 ure 5 (left) shows how the sharpness parameters (γ_e, γ_ℓ) in-
497 fluence readiness behavior and ARS. Larger γ_e effectively
498 suppresses premature answers, while smaller values may
499 over-reward early guesses. Conversely, lower γ_ℓ gently tol-

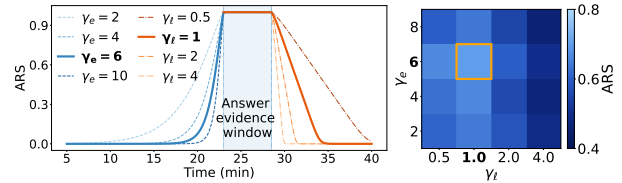


Figure 5. Penalty sharpness for early and late responses. Left: Readiness curves for different penalty strengths. Right: Resulting ARS, with selected γ_e, γ_ℓ combination highlighted.

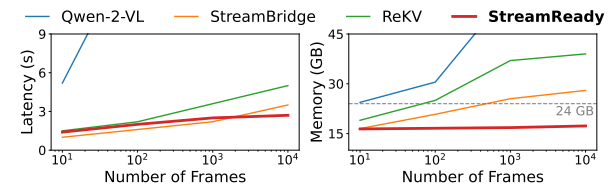


Figure 6. Latency and memory usage analysis.

erates slight delays, whereas higher values penalize hesitation too strongly. Figure 5 (right) shows a broad region of stable performance, with $\gamma_e=6$ and $\gamma_\ell=1$ providing a balanced trade-off between avoiding speculation and allowing realistic delays, demonstrating the robustness of both StreamReady’s readiness behavior and the ARS metric.

Computation Cost and Inference Latency. Figure 6 compares scalability across models as video length increases. Baseline Qwen-2-VL suffers from rapid latency growth and out-of-memory failures due to full-attention accumulation. Retrieval-based (ReKV) and activation-based (StreamBridge) methods improve efficiency but still incur rising overhead: ReKV from expanding KV caches and StreamBridge from repeated costly token-compression steps and heavy activation. In contrast, StreamReady maintains stable latency and memory by using a fixed-size compact memory of centroids and prototypes, keeping retrieval cost minimum. Its readiness mechanism, implemented with a single <RDY> token and lightweight MLP head, adds no extra inference overhead, enabling smooth scalability for long-horizon, readiness-aware streaming.

521 6. Conclusion

522 We present readiness-aware streaming video understand-
523 ing, a formulation that evaluates not only *what* a model an-
524 swers but also *when*, relative to visual evidence. To cap-
525 ture this behavior, we propose the Answer Readiness Score
526 (ARS), a timing-aware metric with asymmetric early and
527 late penalties. Building on this formulation, our proposed
528 framework **StreamReady** integrates long-horizon temporal
529 reasoning with a lightweight readiness mechanism to de-
530 cide when sufficient evidence has appeared. For evaluation,
531 we introduce **ProReady-QA**, a long streaming benchmark
532 with proactive multi-turn questions and annotated answer
533 evidence windows. Together, these contributions establish
534 readiness-aware streaming as a step toward models that an-
535 swer both accurately and on time.

536

References

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

- [1] Marah Abidin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024. 1
- [2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: A visual language model for few-shot learning. In *Advances in Neural Information Processing Systems*, pages 23716–23736, 2022. 1
- [3] Shehreen Azad, Vibhav Vineet, and Yogesh Singh Rawat. Hierarq: Task-aware hierarchical q-former for enhanced video understanding. In *CVPR*, pages 8545–8556, 2025. 1, 2, 3, 6, 7, 8
- [4] Ivana Balazevic, Yuge Shi, Pinelopi Papalampidi, Rahma Chaabouni, Skanda Koppula, and Olivier J Hénaff. Memory consolidation enables long-context video understanding. In *Forty-first International Conference on Machine Learning*, 2024. 2
- [5] Guo Chen, Yin-Dong Zheng, Jiahao Wang, Jilan Xu, Yifei Huang, Junting Pan, Yi Wang, Yali Wang, Yu Qiao, Tong Lu, et al. Videollm: Modeling video sequence with large language models. *arXiv preprint arXiv:2305.13292*, 2023. 1
- [6] Joya Chen, Zhaoyang Lv, Shiwei Wu, Kevin Qinghong Lin, Chenan Song, Difei Gao, Jia-Wei Liu, Ziteng Gao, Dongxing Mao, and Mike Zheng Shou. Videollm-online: Online video large language model for streaming video. In *CVPR*, pages 18407–18418, 2024. 6, 7
- [7] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *Science China Information Sciences*, 67(12):220101, 2024. 6, 7
- [8] Dingxin Cheng, Mingda Li, Jingyu Liu, Yongxin Guo, Bin Jiang, Qingbin Liu, Xi Chen, and Bo Zhao. Enhancing long video understanding via hierarchical event-based memory. *arXiv preprint arXiv:2409.06299*, 2024. 1, 2
- [9] Shangzhe Di, Zhelun Yu, Guanghao Zhang, Haoyuan Li, Tao Zhong, Hao Cheng, Bolin Li, Wanggui He, Fangxun Shu, and Hao Jiang. Streaming video question-answering with in-context video kv-cache retrieval. In *ICLR*, 2025. 1, 2, 7
- [10] Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. In *CVPR*, 2025. 6
- [11] Chaoyou Fu, Haojia Lin, Xiong Wang, Yi-Fan Zhang, Yunhang Shen, Xiaoyu Liu, Haoyu Cao, Zuwei Long, Hetting Gao, Ke Li, et al. Vita-1.5: Towards gpt-4o level real-time vision and speech interaction. *arXiv preprint arXiv:2501.01957*, 2025. 7
- [12] Shenghao Fu, Qize Yang, Yuan-Ming Li, Yi-Xing Peng, Kun-Yu Lin, Xihan Wei, Jian-Fang Hu, Xiaohua Xie, and

Wei-Shi Zheng. Vispeak: Visual instruction feedback in 592

streaming videos. In *ICCV*, 2025. 6, 7 593

[13] Kristen Grauman, Andrew Westbury, Eugene Byrne, 594

Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson 595

Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: 596

Around the world in 3,000 hours of egocentric video. In 597

CVPR, pages 18995–19012, 2022. 5 598

[14] Bo He, Hengduo Li, Young Kyun Jang, Menglin Jia, Xue- 599

fei Cao, Ashish Shah, Abhinav Shrivastava, and Ser-Nam 600

Lim. Ma-Imm: Memory-augmented large multimodal model 601

for long-term video understanding. In *CVPR*, pages 13504– 602

13514, 2024. 1, 2 603

[15] De-An Huang, Shijia Liao, Subhashree Radhakrishnan, 604

Hongxu Yin, Pavlo Molchanov, Zhiding Yu, and Jan Kautz. 605

Lita: Language instructed temporal-localization assistant. In 606

ECCV, pages 202–218. Springer, 2024. 7 607

[16] Qingqiu Huang, Yu Xiong, Anyi Rao, Jiaze Wang, and 608

Dahua Lin. Movienet: A holistic dataset for movie under- 609

standing. In *ECCV*, pages 709–727. Springer, 2020. 5 610

[17] Zhenpeng Huang, Xinhao Li, Jiaqi Li, Jing Wang, Xiangyu 611

Zeng, Cheng Liang, Tao Wu, Xi Chen, Liang Li, and Limin 612

Wang. Online video understanding: Ovbench and videochat- 613

online. In *CVPR*, pages 3328–3338, 2025. 2, 5, 6, 7 614

[18] Minsoo Kim, Kyuhong Shim, Jungwook Choi, and Simyung 615

Chang. Infinipot-v: Memory-constrained kv cache compres- 616

sion for streaming video understanding. In *NeurIPS*, 2025. 617

1, 2, 6, 7 618

[19] Noriyuki Kugo, Xiang Li, Zixin Li, Ashish Gupta, Arpan- 619

deep Khatua, Nidhish Jain, Chaitanya Patel, Yuta Kyuragi, 620

Yasunori Ishii, Masamoto Tanabiki, et al. Videomulti- 621

agents: A multi-agent framework for video question answer- 622

ing. *arXiv preprint arXiv:2504.20091*, 2025. 2 623

[20] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng 624

Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and 625

Chunyuan Li. Llava-onevision: Easy visual task transfer. 626

arXiv preprint arXiv:2408.03326, 2024. 1 627

[21] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, 628

Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Zi- 629

wei Liu, et al. Llava-onevision: Easy visual task transfer. 630

arXiv preprint arXiv:2408.03326, 2024. 6, 7 631

[22] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, 632

Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. 633

Mvbench: A comprehensive multi-modal video understand- 634

ing benchmark. In *CVPR*, pages 22195–22206, 2024. 6 635

[23] Kunchang Li, Xinhao Li, Yi Wang, Yinan He, Yali Wang, 636

Limin Wang, and Yu Qiao. Videomamba: State space model 637

for efficient video understanding. In *European Conference 638**on Computer Vision*, pages 237–255. Springer, 2025. 1 639

[24] Xinhao Li, Yi Wang, Jiashuo Yu, Xiangyu Zeng, Yuhan 640

Zhu, Haiyan Huang, Jianfei Gao, Kunchang Li, Yinan He, 641

Chenting Wang, et al. Videochat-flash: Hierarchical com- 642

pression for long-context video modeling. *arXiv preprint 643**arXiv:2501.00574*, 2024. 2 644

[25] Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An 645

image is worth 2 tokens in large language models. In *ECCV*, 646

pages 323–340. Springer, 2024. 7 647

[26] Yifei Li, Junbo Niu, Ziyang Miao, Chunjiang Ge, Yuanhang 648

Zhou, Qihao He, Xiaoyi Dong, Haodong Duan, Shuangrui 649

- 650 Ding, Rui Qian, et al. Ovo-bench: How far is your video-
651 llms from real-world online video understanding? *arXiv*
652 *preprint arXiv:2501.05510*, 2025. 2, 5, 6
- 653 [27] Junming Lin, Zheng Fang, Chi Chen, Zihao Wan, Fuwen
654 Luo, Peng Li, Yang Liu, and Maosong Sun. Streamingbench:
655 Assessing the gap for mllms to achieve streaming video un-
656 derstanding. *arXiv preprint arXiv:2411.03628*, 2024. 2, 5,
657 6
- 658 [28] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fah-
659 had Khan. Video-ChatGPT: Towards detailed video under-
660 standing via large vision and language models. In *Proceed-*
661 *ings of the 62nd Annual Meeting of the Association for Com-*
662 *putational Linguistics*, 2024. 7
- 663 [29] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra
664 Malik. Egoschema: A diagnostic benchmark for very long-
665 form video language understanding. *Advances in Neural In-*
666 *formation Processing Systems*, 36:46212–46244, 2023. 6
- 667 [30] Tony Montes and Fernando Lozano. Viqagent: Zero-shot
668 video question answering via agent with open-vocabulary
669 grounding validation. *arXiv preprint arXiv:2505.15928*,
670 2025. 2
- 671 [31] Rui Qian, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Shuan-
672 grui Ding, Dahua Lin, and Jiaqi Wang. Streaming long video
673 understanding with large language models. *Advances in Neu-*
674 *ral Information Processing Systems*, 37:119336–119360,
675 2024. 2
- 676 [32] Rui Qian, Shuangrui Ding, Xiaoyi Dong, Pan Zhang, Yuhang
677 Zang, Yuhang Cao, Dahua Lin, and Jiaqi Wang. Dispider:
678 Enabling video llms with active real-time interaction via dis-
679 entangled perception, decision, and reaction. In *CVPR*, 2025.
680 1, 2, 6, 7
- 681 [33] Shuhuai Ren, Linli Yao, Shicheng Li, Xu Sun, and Lu
682 Hou. Timechat: A time-sensitive multimodal large lan-
683 guage model for long video understanding. In *Proceedings*
684 *of the IEEE/CVF Conference on Computer Vision and Pat-*
685 *tern Recognition*, pages 14313–14323, 2024. 7
- 686 [34] Xiaoqian Shen, Yuniang Xiong, Changsheng Zhao, Lemeng
687 Wu, Jun Chen, Chenchen Zhu, Zechun Liu, Fanyu Xiao, Bal-
688 akrishnan Varadarajan, Florian Bordes, et al. Longvu: Spa-
689 tiotemporal adaptive compression for long video-language
690 understanding. *arXiv preprint arXiv:2410.17434*, 2024. 7
- 691 [35] Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng
692 Zhang, Haoyang Zhou, Feiyang Wu, Haozhe Chi, Xun Guo,
693 Tian Ye, Yanting Zhang, et al. Moviechat: From dense token
694 to sparse memory for long video understanding. In *Proceed-*
695 *ings of the IEEE/CVF Conference on Computer Vision and*
696 *Pattern Recognition*, pages 18221–18232, 2024. 7
- 697 [36] Enxin Song, Wenhao Chai, Tian Ye, Jenq-Neng Hwang, Xi
698 Li, and Gaoang Wang. Moviechat+: Question-aware sparse
699 memory for long video question answering. *arXiv preprint*
700 *arXiv:2404.17176*, 2024. 1, 2
- 701 [37] Haibo Wang, Bo Feng, Zhengfeng Lai, Mingze Xu, Shiyu Li,
702 Weifeng Ge, Afshin Dehghan, Meng Cao, and Ping Huang.
703 Streambridge: Turning your offline video large language
704 model into a proactive streaming assistant. In *NeurIPS*, 2025.
705 1, 2, 6, 7, 8
- 706 [38] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan,
707 Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin
Ge, et al. Qwen2-vl: Enhancing vision-language model’s
perception of the world at any resolution. *arXiv preprint*
arXiv:2409.12191, 2024. 1, 6, 7
- [39] Xiaohan Wang, Yuhui Zhang, Orr Zohar, and Serena Yeung-
Levy. Videoagent: Long-form video understanding with
large language model as agent. In *European Conference on*
Computer Vision, pages 58–76. Springer, 2024. 1, 2
- [40] Yueqian Wang, Xiaojun Meng, Yuxuan Wang, Jianxin
Liang, Jiansheng Wei, Huishuai Zhang, and Dongyan Zhao.
Videollm knows when to speak: Enhancing time-sensitive
video comprehension with video-text duet interaction for-
mat. *arXiv preprint arXiv:2411.17991*, 2024. 6
- [41] Yueqian Wang, Xiaojun Meng, Yifan Wang, Huishuai
Zhang, and Dongyan Zhao. Proactivevideoqa: A compre-
hensive benchmark evaluating proactive interactions in video
large language models, 2025. 2, 5
- [42] Yuxuan Wang, Yueqian Wang, Bo Chen, Tong Wu, Dongyan
Zhao, and Zilong Zheng. Omnimmi: A comprehensive
multi-modal interaction benchmark in streaming video con-
texts. In *CVPR*, pages 18925–18935, 2025. 2, 5
- [43] Yuxuan Wang, Cihang Xie, Yang Liu, and Zilong Zheng.
Videollamb: Long-context video understanding with recur-
rent memory bridges. In *ICCV*, 2025. 1, 2, 7
- [44] Ziyang Wang, Shoubin Yu, Elias Stengel-Eskin, Jaehong
Yoon, Feng Cheng, Gedas Bertasius, and Mohit Bansal.
Videotree: Adaptive tree-based video representation for llm
reasoning on long videos. In *Proceedings of the Computer*
Vision and Pattern Recognition Conference, pages 3272–
3283, 2025. 2
- [45] Haomiao Xiong, Zongxin Yang, Jiazuo Yu, Yunzhi Zhuge,
Lu Zhang, Jiawen Zhu, and Huchuan Lu. Streaming video
understanding and multi-round interaction with memory-
enhanced knowledge. *arXiv preprint arXiv:2501.13468*,
2025. 2, 5
- [46] Haolin Yang, Feilong Tang, Linxiao Zhao, Xiang An, Ming
Hu, Huifa Li, Xinlin Zhuang, Boqian Wang, Yifan Lu, Xi-
aofeng Zhang, et al. Streamagent: Towards anticipatory
agents for streaming video understanding. *arXiv preprint*
arXiv:2508.01875, 2025. 7
- [47] Yanlai Yang, Zhuokai Zhao, Satya Narayan Shukla, Aashu
Singh, Shlok Kumar Mishra, Lizhu Zhang, and Mengye Ren.
Streammem: Query-agnostic kv cache memory for stream-
ing video understanding. *arXiv preprint arXiv:2508.15717*,
2025. 1, 2
- [48] Linli Yao, Yicheng Li, Yuancheng Wei, Lei Li, Shuhuai
Ren, Yuanxin Liu, Kun Ouyang, Lean Wang, Shicheng Li,
Sida Li, et al. Timechat-online: 80% visual tokens are
naturally redundant in streaming videos. *arXiv preprint*
arXiv:2504.17343, 2025. 7
- [49] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui,
Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He,
et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv*
preprint arXiv:2408.01800, 2024. 6, 7
- [50] Xiangyu Zeng, Kefan Qiu, Qingyu Zhang, Xinhao Li, Jing
Wang, Jiabin Li, Ziang Yan, Kun Tian, Meng Tian, Xinhai
Zhao, Yi Wang, and Limin Wang. Streamforest: Efficient
online video understanding with persistent event memory,
2025. 2, 5, 6, 7

- 766 [51] Gengyuan Zhang, Tanveer Hannan, Hermine Kleiner, Beste
767 Aydemir, Xinyu Xie, Jian Lan, Thomas Seidl, Volker Tresp,
768 and Jindong Gu. Avila: Asynchronous vision-language agent
769 for streaming multimodal data interaction. *arXiv preprint*
770 *arXiv:2506.18472*, 2025. 2, 8
- 771 [52] Haoji Zhang, Yiqin Wang, Yansong Tang, Yong Liu, Jiashi
772 Feng, Jifeng Dai, and Xiaojie Jin. Flash-vstream: Memory-
773 based real-time understanding for long video streams. *arXiv*
774 *preprint arXiv:2406.08085*, 2024. 1, 7
- 775 [53] Haoji Zhang, Yiqin Wang, Yansong Tang, Yong Liu, Jiashi
776 Feng, Jifeng Dai, and Xiaojie Jin. Flash-vstream: Memory-
777 based real-time understanding for long video streams. In
778 *ICCV*, 2025. 2, 5, 6, 7
- 779 [54] Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng,
780 Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Haoran Tan,
781 Chunyuan Li, and Ziwei Liu. Long context transfer from
782 language to vision. *arXiv preprint arXiv:2406.16852*, 2024.
783 7
- 784 [55] Yuanhan Zhang, Bo Li, haotian Liu, Yong jae Lee, Liangke
785 Gui, Di Fu, Jiashi Feng, Ziwei Liu, and Chunyuan Li. Llava-
786 next: A strong zero-shot video understanding model, 2024.
787 6, 7
- 788 [56] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Zi-
789 wei Liu, and Chunyuan Li. Video instruction tuning with
790 synthetic data. *arXiv preprint arXiv:2410.02713*, 2024. 7
- 791 [57] Zhuo Zhi, Qiangqiang Wu, Wenbo Li, Yinchuan Li, Kun
792 Shao, Kaiwen Zhou, et al. Videoagent2: Enhancing the
793 llm-based agent system for long-form video understanding
794 by uncertainty-aware cot. *arXiv preprint arXiv:2504.04471*,
795 2025. 1, 2
- 796 [58] Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao, Xi
797 Yang, Yongping Xiong, Bo Zhang, Tiejun Huang, and Zheng
798 Liu. Mlvu: A comprehensive benchmark for multi-task
799 long video understanding. *arXiv preprint arXiv:2406.04264*,
800 2024. 6
- 801 [59] Heqing Zou, Tianze Luo, Guiyang Xie, Fengmao Lv,
802 Guangcong Wang, Junyang Chen, Zhuochen Wang, Han-
803 sheng Zhang, Huajian Zhang, et al. From seconds to
804 hours: Reviewing multimodal large language models on
805 comprehensive long video understanding. *arXiv preprint*
806 *arXiv:2409.18938*, 2024. 1