

# ROMA: Real-time Omni-Multimodal Assistant with Interactive Streaming Understanding

Anonymous ACL submission

## Abstract

Recent Omni-multimodal Large Language Models show promise in unified audio, vision, and text modeling. However, streaming audio-video understanding remains challenging, as existing approaches suffer from disjointed capabilities: they typically exhibit incomplete modality support or lack autonomous proactive monitoring. To address this, we present ROMA, a **real-time omni-multimodal assistant for unified reactive and proactive interaction**. ROMA processes continuous inputs as synchronized *multimodal units*, aligning dense audio with discrete video frames to handle granularity mismatches. For online decision-making, we introduce a lightweight *speak head* that decouples response initiation from generation to ensure precise triggering without task conflict. We train ROMA with a curated streaming dataset and a two-stage curriculum that progressively optimizes for streaming format adaptation and proactive responsiveness. To standardize the fragmented evaluation landscape, we reorganize diverse benchmarks into a unified suite covering both proactive (alert, narration) and reactive (QA) settings. Extensive experiments across 12 benchmarks demonstrate ROMA achieves state-of-the-art performance on proactive tasks while competitive in reactive settings, validating its robustness in unified real-time omni-multimodal understanding. Code and benchmark are available [here<sup>1</sup>](#).

## 1 Introduction

Recent advances in omni-multimodal large language models (OLLMs), such as GPT-4o (Hurst et al., 2024), have enabled unified modeling of speech, vision, and text. This progress facilitates real-world streaming audio-video understanding, defined as combining **reactive** and **proactive** capabilities (Figure 1). In the reactive setting, the model

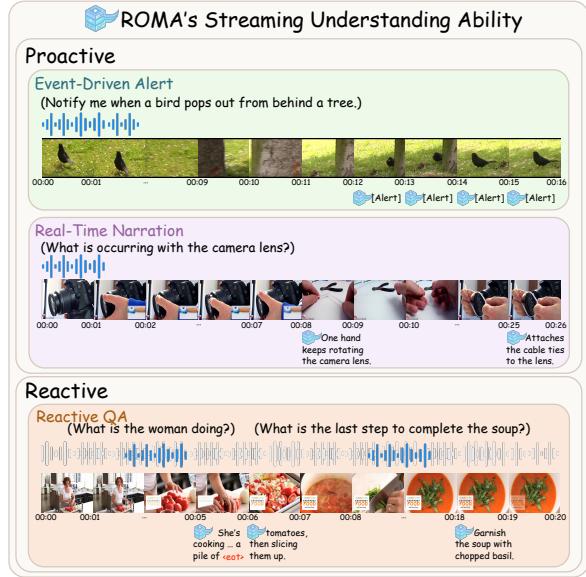


Figure 1: **ROMA’s streaming understanding capabilities.** It supports proactive tasks, including event alerts and narration, alongside reactive question answering.

answers after the query, whereas in the proactive setting, it follows an instruction to continuously monitor the input stream and respond only when conditions are met. Unifying these capabilities is vital for real-world utility, yet the divergent interaction paradigms make it challenging (Horvitz, 1999; Xi et al., 2025; Driess et al., 2023).

Despite the critical need for such unification, existing studies typically lack unified modality support and streaming capabilities. Specifically, speech-centric streaming models (Défossez et al., 2024; Zhang et al., 2025) focus on audio generation but lack visual perception. Conversely, while some approaches address streaming video understanding (Chen et al., 2024; Zhang et al., 2024b), they typically neglect synchronized audio and are confined to specific tasks (e.g., alert or narration). Consequently, unified streaming audio-video understanding remains largely under-explored.

To realize such unification faces two challenges.

<sup>1</sup><https://anonymous.4open.science/r/ROMA-6746>

First, audio and video exhibit mismatched temporal granularities. While naturally synchronized, audio signals are dense and continuous, whereas video comprises sparse, discrete frames. Under such heterogeneity, maintaining robust cross-modal alignment and fusion demands precise synchronization. Second, effective streaming interaction requires real-time proactive decision-making. Upon integrating these asynchronous signals, the model must continuously synthesize context to determine both response timing and content, conditioned strictly on the stream prefix.

To address these challenges, we propose ROMA, a **Real-time Omni-Multimodal Assistant** with interactive streaming understanding. To tackle the granularity mismatch, ROMA segments continuous audio into one-second intervals synchronized with video frames, forming temporally aligned units that are processed sequentially as the stream unfolds. We further adapt chunked Time-aligned Multimodal RoPE (TMRoPE) (Xu et al., 2025a) to enforce a shared temporal timeline. For proactive decision-making, ROMA introduces a lightweight speak head parallel to the standard language modeling (LM) head to explicitly predict response timing, decoupling timing from content generation to prevent task interference. Finally, we support this system with a custom streaming dataset and a two-stage training curriculum, progressively optimizing the model for cross-modal streaming format adaptation and proactive responsiveness.

For a comprehensive evaluation, streaming audio-video understanding demands assessing both reactive and proactive capabilities. However, as compared in Table 1, existing benchmarks suffer from inconsistent taxonomies and fragmented protocols, often failing to cover both interaction modes. To enable unified comparison, we reorganize the evaluation landscape into two standardized settings: a *proactive* mode that tests the ability to autonomously trigger responses at precise moments, and a *reactive* mode that emphasizes understanding temporal evolution in standard QA. Empirically, ROMA consistently outperforms existing streaming VideoLLMs across both modes. Furthermore, evaluations on open-ended audio-query QA against open-source OLLMs confirm its superior capability in unified audio-video understanding.

In summary, our contributions are as follows:

- **Unified streaming framework:** We formally define the task of streaming audio-video un-

derstanding and propose ROMA, an omni-multimodal assistant unifying reactive and proactive capabilities, supported by a curated dataset and a two-stage curriculum.

- **Standardized evaluation benchmark:** We establish a comprehensive streaming benchmark by reorganizing fragmented tasks into unified *reactive* and *proactive* settings to facilitate rigorous and consistent comparison.
- **Superior performance and analysis:** ROMA achieves state-of-the-art results across proactive benchmarks while competitive on reactive and open-ended QA. Extensive analysis verifies the efficacy of our timing mechanisms and training strategies.

Benchmark	Alert	Narration	Reactive QA
StreamingBench (Lin et al., 2024)	✓	✗	✓
StreamBench (Wu et al., 2024a)	✗	✗	✓
OVO-Bench (Niu et al., 2025)	✗	✓	✓
SVBench (Yang et al., 2025b)	✗	✗	✓
OmniMMI (Wang et al., 2025b)	✓	✗	✓
OVBench (Huang et al., 2025)	✗	✗	✓
Ours	✓	✓	✓

Table 1: Coverage of key streaming ability across representative streaming video benchmarks.

## 2 Related Works

**Reactive Models** Most existing streaming systems are studied in the reactive setting, answering only after the query arrives. Within this regime, memory-based methods maintain long-range context for coherent understanding over evolving streams (Qian et al., 2024; Zhang et al., 2024a; Wang et al., 2024b; Zhang et al., 2024b; Xiong et al., 2025; Wang et al., 2025a; Zhao et al., 2025), and KV-cache based methods optimize efficiency via scheduling or compression (Di et al., 2025; Ning et al., 2025; Yang et al., 2025a; Xu et al., 2025c; Chen et al., 2025b). Recent omni-multimodal models also adhere to this reactive protocol: MiniCPM-o 2.6 (Yao et al., 2024), Qwen2.5-Omni (Xu et al., 2025a), and Qwen3-Omni (Xu et al., 2025b) support low-latency interaction, and Stream-Omni (Zhang et al., 2025) enables visually-conditioned speech generation, yet none explicitly model proactive monitoring and triggering.

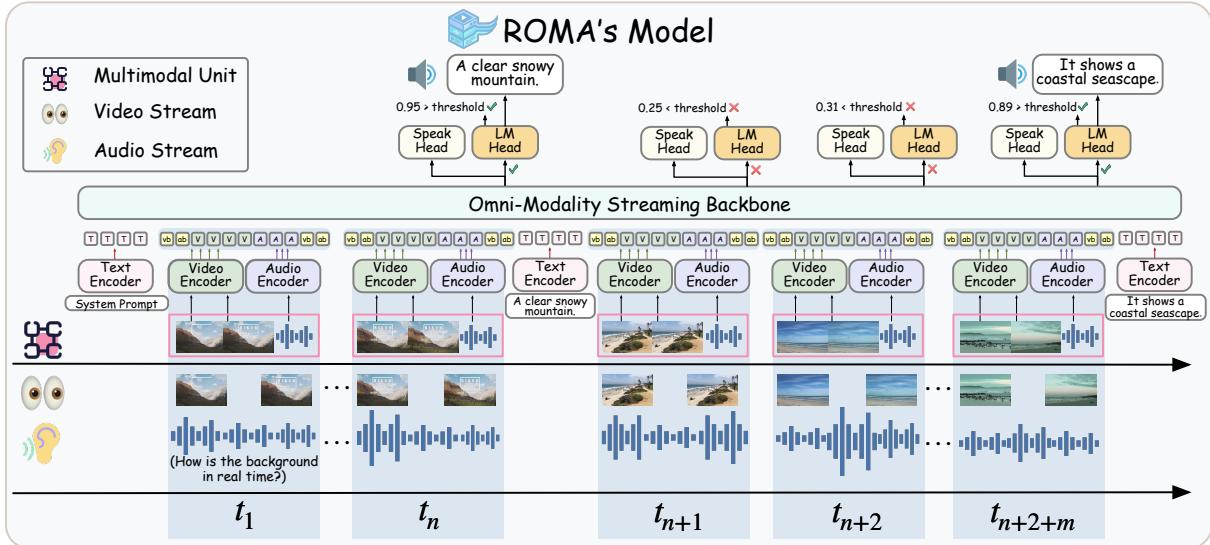


Figure 2: **Model Architecture.** Streaming inputs are processed as aligned multimodal units. The speak head determines response timing, activating the LM head (illustrated via narration) upon crossing a probability threshold.

**Proactive Models** In contrast, proactive streaming prioritizes continuous monitoring and time-sensitive triggering (e.g., alerts and real-time narration). Proactive VideoLLMs leverage online formats or explicit decision modeling to determine intervention timing (Chen et al., 2024; Yang et al., 2025d; Li et al., 2025; Yang et al., 2025c; Qian et al., 2025), with some explicitly targeting live narration (Chen et al., 2025a). However, these approaches remain predominantly video-centric, neglecting streaming audio.

**Streaming Video Understanding Benchmarks** Recent benchmarks prioritize time-sensitive, interactive evaluation (Table 1). StreamingBench (Lin et al., 2024) and OVO-Bench (Niu et al., 2025) assess temporal perception, while StreamBench (Wu et al., 2024a) and SVBench (Yang et al., 2025b) focus on long-horizon memory. Moreover, OmniIMMI (Wang et al., 2025b) and OVBBench (Huang et al., 2025) incorporate proactive capabilities, including real-time narration and alerts.

**Tables 1 and 10 summarizes prior works.**

### 3 Method

To unify reactive answering and proactive timing over continuous inputs, ROMA integrates architectural designs with a tailored training strategy. Section 3.1 presents the model architecture, utilizing chunked TMRoPE and a speak head for alignment and timing control. Section 3.2 details the training and inference pipeline, encompassing dataset construction and a two-stage fine-tuning recipe.

### 3.1 Model Architecture

As illustrated in Figure 2, ROMA processes streaming omni-modal inputs via a unified LLM backbone. We introduce a speak head parallel to the LM head to decouple interaction timing from content generation. This architecture addresses temporal alignment and proactive decision-making through the following mechanisms.

**Multimodal units for temporally aligned streaming inputs.** To support unified streaming understanding across modalities, we organize audio and video into fixed-interval multimodal units. Following the input format and tokenization of Qwen2.5-Omni, we treat all audio and video signals within each one-second interval as a unit. We align audio with video frames sampled from the same interval, extract their features, and wrap them with special tokens. This retains Qwen2.5-Omni’s native format for compatibility while grounding audio in the preceding visual context:

```
<|vision_bos|><|audio_bos|> [video tokens] [audio tokens] <|audio_eos|><|vision_eos|>
```

These multimodal units are fed into the LLM backbone sequentially as the stream unfolds. This process ensures that the model continuously accumulates aligned cross-modal context from the stream prefix, establishing a temporal basis for subsequent causal decision-making.

**Chunk-Level Temporal Position Encoding** We adapt Qwen2.5-Omni’s Time-aligned Multimodal RoPE (TMRoPE) to chunked audio-video streams to support incremental encoding as units arrive.

209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223

Each one-second unit interleaves visual and auditory tokens, assigning time-aligned 3D position IDs (temporal, height, and width) to preserve their cross-modal correspondence. Consistent with the pre-trained vision encoder, multi-frame visual inputs are temporally aggregated into a fused representation during encoding. All video tokens within a unit therefore share a constant temporal ID. In contrast, audio tokens retain fine-grained temporal IDs at a 40ms resolution to preserve auditory temporal fidelity. To ensure boundary alignment, `<|vision_bos|>` and `<|audio_bos|>` share the same base position ID, subsequent units extend the global timeline by continuing from the maximum position ID of the previous unit.

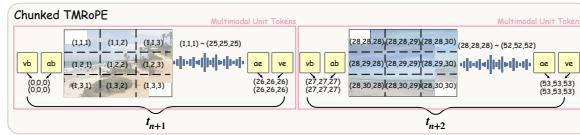


Figure 3: **Chunked TMRoPE**. Seamlessly extends the global timeline to streaming inputs by assigning cumulative positional IDs across discrete units.

224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242

**Speak Head** To enable autonomous intervention timing, we design a lightweight speak head. As illustrated in Figure 2, this module is implemented as a two-layer MLP, parallel to the LM head, on top of the streaming backbone. Upon processing each multimodal unit (one second of context), the speak head evaluates the current stream prefix and outputs a probability for a binary decision indicating whether a response is required. A response is triggered if this probability exceeds a threshold; otherwise, the model remains silent and continues consuming the stream. This design decouples the timing decision from text generation, mitigating interference from generative biases. Leveraging findings that upper layers encode high-level features (Tenney et al., 2019; Belrose et al., 2023), we compute the speak head input as a learnable weighted combination of hidden states from the last  $K$  layers, with  $K=4$  in our experiments.

## 243 3.2 Training and Inference Pipeline

### 244 3.2.1 Dataset Construction

245 To enable end-to-end proactive and reactive supervision, we construct a comprehensive streaming dataset structured into two categories and three sub-tasks (Figure 4). Detailed processing pipelines  
246 are provided in Appendix A.4.

250 To equip the model with  
251 the ability to continuously monitor streams and trigger alerts, we curate data from DiDeMo (Anne Hendricks et al., 2017), OOPS (Epstein et al., 2020),  
252 and Charades-STA (Zhou et al., 2018). We reformulate these samples into alert-style tasks (e.g.,  
253 “Alert me when [event] happens”) to train the model  
254 in event-driven temporal grounding.

255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277

**Online Narration (109K)** To foster continuous event tracking and incremental summarization, we construct narration samples from MM-DuetIT (Wang et al., 2024a), COIN (Tang et al., 2019), YouCook2 (Zhou et al., 2018), and ActivityNet (Caba Heilbron et al., 2015). Unlike prior works that use dense supervision, we specifically train the model to generate captions only at segment transitions, enabling it to provide concise, real-time updates as the visual context evolves.

**Reactive QA (540K)** To stabilize general audio-video understanding, we aggregate large-scale reactive QA data from InternVid (Wang et al., 2023), CogStream (Zhao et al., 2025), and others (Chen et al., 2023; Yang et al., 2022; Yao et al., 2025; Fu et al., 2025). These samples cover past events, temporal ordering, and future reasoning.

To ensure unified processing, we synthesize text queries into speech, training the model to handle audio instructions under streaming inputs.

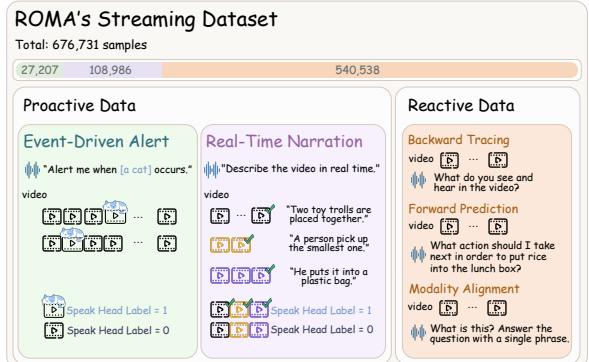


Figure 4: **Overview of ROMA’s Streaming Dataset**. Statistics, task taxonomy, and sample formats.

### 278 3.2.2 A Two-Stage Fine-Tuning Recipe

279 Training an end-to-end streaming omni-multimodal  
280 model from scratch is computationally prohibitive.  
281 We fundamentally view streaming capability as  
282 a transfer problem: adapting a strong foundation  
283 model optimized for processing complete videos to  
284 handle incremental streams. We thus propose a simple  
285 yet effective two-stage recipe. Stage 1 adapts

286 the model to the streaming multimodal input for-  
287 mat, while Stage 2 learns precise response timing  
288 and proactive policies. In both stages, we freeze all  
289 encoders and fine-tune the remaining parameters  $\theta$ .  
290

**Stage 1: Streaming Template Alignment** This  
291 stage mitigates the distribution shift between of-  
292 fline training and streaming inference. We utilize  
293 reactive QA datasets to adapt the model to the mul-  
294 timodal unit streaming format. Samples are restructured  
295 into sequential units  $X$  to simulate streaming,  
296 with the audio query and text response  $Y$  appended.

We optimize the standard autoregressive lan-  
guage modeling objective over the response tokens.  
Let  $\mathcal{D}_{\text{QA}}$  denote the reactive QA dataset. For a  
sample  $(X, Y) \sim \mathcal{D}_{\text{QA}}$ , where  $Y = \{y_1, \dots, y_L\}$   
represents the answer sequence, the loss is:

$$\mathcal{L}_{\text{LM}} = -\mathbb{E}_{(X, Y) \sim \mathcal{D}_{\text{QA}}} \left[ \sum_{i=1}^L \log P(y_i | y_{<i}, X; \theta) \right].$$

This stage ensures the model retains its audio-video  
understanding while adapting to streaming inputs.

**Stage 2: Time-Aware Decision Making** With  
the backbone adapted to streaming inputs, this  
stage activates the speak head to learn *when* to  
respond. We formulate response timing as a bi-  
nary classification task at each multimodal unit  
step. The positive labels are task-dependent: for  
proactive alerts, valid triggers lie within the event  
window; for narration, they align with segment  
boundaries. To mitigate trigger sparsity, we bal-  
ance the loss using  $w_{\text{pos}} = N_{\text{neg}}/N_{\text{pos}}$  derived  
from dataset statistics.

Let  $p_t$  be the speak head’s predicted probability  
at time step  $t$ , and  $z_t \in \{0, 1\}$  be the ground truth  
label. The timing loss is formulated as a weighted  
Binary Cross-Entropy (BCE):

$$\begin{aligned} \mathcal{L}_{\text{time}} = -\mathbb{E}_{X \sim \mathcal{D}_{\text{stream}}} & \left[ \frac{1}{T} \sum_{t=1}^T \left( w_{\text{pos}} z_t \log p_t \right. \right. \\ & \left. \left. + (1 - z_t) \log(1 - p_t) \right) \right]. \end{aligned}$$

To prevent generation quality degradation while  
optimizing purely for timing, we mix a small por-  
tion of the Stage 1 reactive QA data ( $\mathcal{D}_{\text{QA}}$ ) during  
training. The final objective is a joint optimization:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{time}} + \lambda \cdot \mathcal{L}_{\text{LM}},$$

where  $\mathcal{L}_{\text{LM}}$  is calculated only on the mixed QA  
samples to maintain linguistic competence, and  $\lambda$   
balances the two objectives.

### 3.2.3 Inference Procedure

During inference, we strictly follow the training  
configuration. Video frames are uniformly sam-  
pled at 2 fps, and each frame is resized so that  
the number of pixels does not exceed 65,536. We  
maintain a persistent KV cache across the stream,  
so each step only encodes the current multimodal  
unit. Under this setup, encoding one unit takes  
0.3697 seconds on average.

## 4 Unified Streaming Evaluation Framework

Effective streaming understanding demands mod-  
els capable of answering queries and autonomously  
determining interaction timing. Addressing the  
fragmentation in existing benchmarks (Table 1),  
we establish a unified framework comprising two  
primary settings: **proactive interaction**, where the  
model autonomously monitors the stream to trig-  
ger responses, and **reactive interaction**, where it  
answers queries based on accumulated context.

### 4.1 Proactive Streaming Interaction

In the proactive setting, the model receives an in-  
struction at the start and must process the stream  
to determine both the precise timing and content of  
the response. We categorize this into two sub-tasks:  
event-driven alert and real-time narration.

#### 4.1.1 Event-Driven Alert

This task evaluates the model’s temporal awareness,  
specifically its ability to detect transient events and  
trigger immediate notifications. We assess this ca-  
pability under two settings.

**Static Temporal Grounding.** Following MM-  
Duet on QVHighlights (Lei et al., 2021) and  
Charades-STA (Gao et al., 2017), ROMA incre-  
mentally predicts response probabilities for each  
multimodal unit. For QVHighlights, we rank times-  
tamps by normalized probabilities, reporting mAP  
(ranking quality) and HIT@1 (top-1 accuracy). For  
localization on Charades-STA, we threshold prob-  
abilities to predict spans, reporting R@0.5 and  
R@0.7 (recall at 0.5 and 0.7 temporal overlap).

**Dynamic Streaming Decision.** This configura-  
tion enforces a strict streaming protocol where the  
model makes instantaneous decisions conditioned  
exclusively on the current multimodal unit. We  
conduct a comprehensive evaluation across OmniMMI  
(PA), StreamingBench (PO), and OVO-  
Bench (CRR, REC), spanning both single-event

367 alerts and multi-event recurrence. Specifically,  
 368 for OVO-Bench, we reformulate the original QA-  
 369 centric annotations into streaming alert targets to  
 370 evaluate instantaneous responsiveness. To mitigate  
 371 transient probability fluctuations, we employ a slid-  
 372 ing window mechanism. Success is determined by  
 373 the temporal inclusion of the autonomously trig-  
 374 gered response within the ground-truth interval.  
 375

See Appendix A.3 for detailed settings.

#### 376 4.1.2 Real-Time Narration

377 We define streaming narration as the incremental  
 378 summarization of evolving events devoid of future  
 379 context. To evaluate this capability, we employ  
 380 two settings: a continuous YouCook2 adaptation,  
 381 constructed by concatenating annotated clips to en-  
 382 force generation at segment transitions, and the  
 383 OVO-Bench (SSR) task, where responses are trig-  
 384 gered via prediction thresholds and appended to the  
 385 streaming context. Performance is assessed using  
 386 the F1 score for temporal localization, BERTScore  
 387 for the semantic quality of aligned responses, and a  
 388 GPT-4o-based evaluation of coherence, alignment,  
 389 and conciseness (detailed in Figure 12).

#### 390 4.2 Reactive QA

391 In the reactive setting, the model must interpret  
 392 temporal evolution to answer questions constrained  
 393 to the causal video history. We utilize OVO-Bench  
 394 and StreamingBench for standardized evaluation,  
 395 employing text-based queries to ensure fairness  
 396 against VideoLLMs baselines and reporting accu-  
 397 racy. To further approximate real-world interac-  
 398 tion, we extend the assessment to Video-MME and  
 399 EgoSchema using synthesized speech inputs. This  
 400 setting evaluates comprehensive audio-video un-  
 401 derstanding, with open-ended responses scored by  
 402 GPT-4o (detailed in Figure 11).

### 403 5 Experiment

#### 404 5.1 Implementation Details

405 To address trigger sparsity, we set the positive  
 406 weight  $w_{pos} = 3$  in the weighted BCE loss. For  
 407 inference, we adopt a pipelined real-time approxi-  
 408 mation: the model processes unit  $t$  while simulta-  
 409 neously acquiring unit  $t + 1$ . To ensure synchro-  
 410 nization, we cap generation at 25 tokens (approx.  
 411 1s) per segment, allowing longer responses to con-  
 412 tinue across subsequent units. Please refer to Ap-  
 413 pendix A.5 for detailed training configurations and  
 414 complete decoding protocols.

## 5.2 Experimental Results

**Baseline Methods** In the proactive setting, we limit comparison to streaming-capable models: VideoLLM-Online (the basis for many efficiency-focused architectures), MMDuet, and Dispider. We reproduce results using accessible implementations, defaulting to reported figures otherwise. For Reactive QA, we benchmark against representative streaming VideoLLMs. To assess full-modality understanding, we extend evaluation to open-source omni-modal models, including Qwen2.5-Omni, MiniCPM-o, and VITA-1.5.

Method	QVHighlight Charades-STA	
	mAP / HIT@1	R@0.5 / 0.7
TimeChat	14.5 / 23.9	32.2 / 13.4
VTimeLLM	–	31.2 / 11.4
HawkEye	–	31.4 / 14.5
VTG-LLM	16.5 / 33.5	33.8 / 15.7
MMDuet	31.3 / 49.6	42.4 / 18.0
<b>Ours</b>	<b>53.7 / 53.0</b>	<b>44.3 / 19.9</b>
<i>- Ablation Study</i>		
Mixed Training	50.3 / 44.7	28.2 / 10.1
$K = 1$	46.4 / 47.4	32.4 / 13.1
<i>- Sensitivity Analysis</i>		
$w_{pos} = 2$	47.5 / 52.5	42.2 / 18.4
$w_{pos} = 4$	47.3 / 49.1	38.0 / 16.4

Table 2: Comparison with existing methods on QVHighlights and Charades-STA benchmarks.

Method	PA	PO	CRR	REC
VideoLLM-online	0.50	4.13	27.08	14.29
MMDuet	22.00	29.44	16.67	12.77
Dispider	–	25.34	<b>48.75</b>	18.05
M4-a	25.50	–	–	–
<b>Ours</b>	<b>37.50</b>	<b>53.60</b>	35.42	<b>33.81</b>
<i>- Ablation Study</i>				
Mixed Training	34.50	50.80	25.00	13.13
w/o Speak Head	12.50	12.00	0.00	6.46
$K = 1$	26.00	56.40	31.25	24.32
<i>- Sensitivity Analysis</i>				
$w_{pos} = 2$	31.00	52.76	39.58	31.54
$w_{pos} = 4$	31.00	52.15	37.50	26.74

Table 3: Comparison across single-alert (PA, PO, CRR) and recurring alert (REC) benchmarks.

**Event-Driven Alert** In static temporal grounding (Table 2), ROMA advances temporal localization on QVHighlights (53.7 mAP) and Charades-STA

(44.3/19.9 R@0.5/0.7), confirming that incremental speak probabilities provide enhanced temporal saliency for precise ranking and prediction. In the dynamic setting (Table 3), ROMA demonstrates strong efficacy on single-alert tasks: it excels on PA and PO while remaining competitive on CRR, validating its precise proactive triggering and robust evidence accumulation. Furthermore, ROMA dominates on the REC benchmark, validating its recurrence modeling for tracking repeated instances.

**Real-Time Narration** As shown in Table 4, ROMA achieves the best temporal triggering accuracy, obtaining an F1 score of 35.21 on YouCook2 and 14.54 on OVO-Bench (SSR), which indicates more precise alignment between generated responses and the annotated narration windows. It also achieves the highest GPT-4o score on both benchmarks. This score averages three criteria (story coherence, alignment to ground truth, and conciseness), with the per-criterion breakdown in Table 9, suggesting more coherent and better-aligned narration when generation is triggered online and the outputs are carried forward as context.

Method	YouCook2			OVO-Bench (SSR)		
	F1	BERT	GPT	F1	BERT	GPT
TimeChat	21.70	–	–	–	–	–
VTG-LLM	17.50	–	–	–	–	–
VideoLLM-online	18.82	0.82	0.17	10.24	<b>0.84</b>	0.18
MMDuet	17.81	0.83	0.23	9.02	0.79	0.31
<b>Ours</b>	<b>35.21</b>	<b>0.83</b>	<b>0.39</b>	<b>14.54</b>	0.83	<b>0.42</b>
<i>- Ablation Study</i>						
Mixed Training	31.42	0.81	0.34	8.88	0.80	0.33
w/o Speak Head	9.25	0.79	0.24	3.39	0.77	0.26
$K = 1$	34.43	0.82	0.37	9.64	0.78	0.32
<i>- Sensitivity Analysis</i>						
$w_{pos} = 2$	27.82	0.83	0.45	10.38	0.57	0.38
$w_{pos} = 4$	35.55	0.81	0.47	13.48	0.75	0.34

Table 4: Streaming narration results on YouCook2 and OVO-Bench (SSR). We report F1 for temporal window alignment, and use BERTScore and averaged GPT-4o scores to assess narration quality.

**Reactive QA** On OVO-Bench (Table 6), ROMA leads in both “Real-time Visual Perception” and “Backward Tracing”. Its superiority over streaming baselines highlights enhanced sensitivity to time-localized cues and robust utilization of historical

evidence under truncated contexts. On StreamingBench (Table 7), ROMA maintains high accuracy and secures the top rank on “Omni-Source Understanding” benchmark, attributed to preserving aligned audio during training, which bolsters audio–visual integration. In full-modality evaluation (Table 5), ROMA attains the best performance on Video-MME (without subtitles) and remains competitive on EgoSchema. Notably, these results utilize spoken queries with joint audio–visual inputs to approximate conversational interaction, distinct from text-prompted prior work.

Overall, ROMA strengthens temporal awareness and streaming decision-making, optimizing timing and content via audio–video joint modeling.

Method	Video-MME	EgoSchema
Qwen2.5-Omni	20.50	<b>58.40</b>
VITA-1.5	28.56	45.40
MiniCPM-o	19.37	55.20
<b>Ours</b>	<b>33.30</b>	55.40
<i>- Ablation Study</i>		
Mixed Training	33.00	50.20
w/o speak head	9.11	12.80
$K = 1$	34.56	54.00
<i>- Sensitivity Analysis</i>		
$w_{pos} = 2$	33.20	52.60
$w_{pos} = 4$	33.10	54.80

Table 5: Full-modality QA results on Video-MME (no subtitles) and EgoSchema, evaluated with spoken questions to approximate real conversational interaction.

### 5.3 Ablation Study

**Single-Stage vs. Two-Stage Training** We validate the two-stage curriculum by mixing all data and training directly with the stage-2 objective. This variant consistently degrades on tasks that require online timing and triggering, most notably on dynamic decision making (e.g., REC) and streaming narration (Table 3, Table 4). The results indicate that progressive training is important for learning well-calibrated temporal decision making under streaming input.

**Speak Head for Response Gating** We replace the speak head with a ‘<silence>’ token following prior work, and cast triggering as next-token prediction with a reweighted loss. Lacking explicit probabilities, we omit QVHighlights and Charades-STA, instead evaluating triggering based on the first non-‘<silence>’ token.

Method	Real-time Visual Perception						Backward Tracing		
	OCR	ACR	ATR	STU	FPD	OJR	EPM	ASI	HLD
VideoLLM-online	8.05	23.85	12.07	14.04	45.54	21.20	22.22	18.80	12.18
MMDuet	13.42	11.93	14.66	11.80	14.85	10.33	10.44	8.78	0.54
Dispider	57.72	49.54	62.07	<b>44.94</b>	61.39	51.63	48.48	<b>55.41</b>	4.30
Flash-VStream-7B	24.16	29.36	28.45	33.71	25.74	28.80	39.06	37.16	5.91
<b>Ours</b>	<b>63.09</b>	<b>53.21</b>	<b>68.10</b>	39.33	<b>69.31</b>	<b>58.15</b>	<b>55.89</b>	47.30	<b>23.66</b>
<i>- Ablation Study</i>									
Mixed Training	63.09	55.05	63.79	37.64	61.39	55.43	55.22	45.95	27.96
w/o Speak Head	61.07	55.05	63.97	39.89	65.35	54.89	53.87	47.97	29.03
$K = 1$	61.47	55.05	68.10	39.89	65.35	60.33	56.57	46.62	20.97
<i>- Sensitivity Analysis</i>									
$w_{pos} = 2$	64.43	51.38	68.97	39.33	64.36	60.87	54.88	46.62	20.97
$w_{pos} = 4$	65.10	54.13	68.97	38.20	70.30	61.41	56.57	46.27	22.58

Table 6: Reactive QA results on OVO-Bench (excluding Forward Active Responding), evaluating time-sensitive understanding across Real-time Visual Perception and Backward Tracing.

Method	Real-Time Visual Understanding										Omni-Source Understanding				Contextual Understanding		
	OP	CR	CS	ATP	EU	TR	PR	SU	ACP	CT	ER	SCU	SD	MA	ACU	MCU	SQA
VideoLLM-Online	39.07	40.06	34.49	31.05	45.96	32.40	31.48	34.16	42.49	27.89	31.20	26.51	24.10	32.00	24.19	29.20	26.55
Flash-VStream	25.89	43.57	24.91	23.87	27.33	13.08	18.52	25.20	23.87	48.70	25.91	24.90	25.60	28.40	24.80	25.20	24.12
Dispider	74.92	75.53	74.10	73.08	74.44	59.52	76.14	<b>62.91</b>	62.16	45.80	35.46	25.26	38.57	43.34	<b>39.62</b>	27.65	33.61
<b>Ours</b>	<b>76.96</b>	<b>78.91</b>	<b>77.92</b>	<b>82.05</b>	<b>74.84</b>	<b>72.90</b>	<b>82.41</b>	61.79	<b>65.91</b>	<b>51.06</b>	<b>40.40</b>	<b>34.80</b>	<b>50.40</b>	<b>58.80</b>	37.60	<b>34.00</b>	<b>44.47</b>
<i>- Ablation Study</i>																	
Mixed Training	75.51	85.71	76.19	78.23	59.77	61.05	73.21	60.00	59.67	23.38	38.80	26.00	40.80	47.20	35.60	27.20	26.80
w/o Speak Head	76.13	70.49	74.14	82.40	72.86	70.80	84.78	63.20	64.91	51.69	39.75	30.36	45.87	47.20	35.27	24.79	24.80
$K = 1$	76.69	82.03	78.86	82.05	74.84	72.90	79.63	59.76	64.49	50.53	38.80	29.60	46.40	51.60	33.20	27.60	22.80
<i>- Sensitivity Analysis</i>																	
$w_{pos} = 2$	75.61	81.25	76.97	82.37	71.70	75.08	81.48	62.20	65.62	50.00	39.60	30.80	59.38	52.40	35.60	28.00	26.40
$w_{pos} = 4$	75.61	80.47	79.18	82.37	73.58	75.08	82.41	63.10	65.34	47.34	39.60	28.80	44.40	51.60	35.60	28.40	26.40

Table 7: Reactive QA results on StreamingBench (excluding PO), evaluating real-time understanding under streaming input across Real-Time Visual Understanding, Omni-Source Understanding, and Contextual Understanding.

**Last-Layer vs. Last-4-Layer Aggregation** We ablate four-layer aggregation by restricting the speak head to the final layer ( $K=1$ ). This notably degrades temporal grounding and dynamic triggering (Tables 2, 3) while leaving timestamp-conditioned understanding largely unaffected (Tables 6, 7). This confirms multi-layer aggregation yields robust signals essential for streaming.

#### 5.4 Sensitivity analysis

We sweep the positive weight  $w_{pos}$  in the weighted BCE loss of the speak head to mitigate the class imbalance from sparse speaking timestamps. We observe that  $w_{pos}$  is critical for proactive tasks (Tables 2–4), while reactive understanding and full-modality QA remain insensitive (Tables 5, 6). Overall,  $w_{pos} = 3$  yields the most balanced perfor-

mance. See Appendix A.2 for sensitivity analysis on inference-time triggering thresholds.

## 6 Conclusion

We introduce ROMA, a real-time omni-multimodal assistant that redefines streaming interaction as the unification of proactive and reactive paradigms. ROMA is the first framework to excel in both modes. To achieve this, we construct a streaming dataset and training recipe that enhance temporal modeling and decision-making. Furthermore, we standardize evaluation through a unified protocol tailored to this dual paradigm, where ROMA demonstrates superior performance. Finally, we provide a systematized overview of prior methods to facilitate future research.

## 523 Limitations

524 While optimized for streaming interaction, the  
525 model remains susceptible to distortions such as  
526 signal degradation and audio-video asynchrony.  
527 Additionally, while capable of continuous stream-  
528 ing, capturing extremely long-term dependencies  
529 spanning hours remains constrained by finite con-  
530 text windows and memory. Finally, optimizing the  
531 trade-off between inference efficiency and response  
532 quality under strict resource constraints remains a  
533 critical direction for future work.

## 534 Ethical Statement

535 This work utilizes publicly available datasets con-  
536 sistent with their original licenses. While ROMA  
537 enables proactive monitoring capabilities, we ac-  
538 knowledge the potential risk of misuse for unau-  
539 thorized surveillance or privacy infringement. This  
540 model is intended for research purposes; due to  
541 the possibility of hallucinations or biases inherited  
542 from the base LLM, human oversight is strictly  
543 required for critical real-world applications.

## 544 References

545 Lisa Anne Hendricks, Oliver Wang, Eli Shechtman,  
546 Josef Sivic, Trevor Darrell, and Bryan Russell. 2017.  
547 Localizing moments in video with natural language.  
548 In *Proceedings of the IEEE international conference  
549 on computer vision*, pages 5803–5812.

550 Nora Belrose, Zach Furman, Logan Smith, Danny Ha-  
551 lawi, Igor Ostrovsky, Lev McKinney, Stella Bider-  
552 man, and Jacob Steinhardt. 2023. Eliciting latent  
553 predictions from transformers with the tuned lens.  
554 *arXiv preprint arXiv:2303.08112*.

555 Fabian Caba Heilbron, Victor Escorcia, Bernard  
556 Ghanem, and Juan Carlos Niebles. 2015. Activitynet:  
557 A large-scale video benchmark for human activity  
558 understanding. In *Proceedings of the ieee conference  
559 on computer vision and pattern recognition*, pages  
560 961–970.

561 Joya Chen, Zhaoyang Lv, Shiwei Wu, Kevin Qinghong  
562 Lin, Chenan Song, Difei Gao, Jia-Wei Liu, Ziteng  
563 Gao, Dongxing Mao, and Mike Zheng Shou. 2024.  
564 Videollm-online: Online video large language  
565 model for streaming video. In *Proceedings of the  
566 IEEE/CVF Conference on Computer Vision and Pat-  
567 tern Recognition*, pages 18407–18418.

568 Joya Chen, Ziyun Zeng, Yiqi Lin, Wei Li, Zejun Ma,  
569 and Mike Zheng Shou. 2025a. Livecc: Learning  
570 video llm with streaming speech transcription at scale.  
571 In *Proceedings of the Computer Vision and Pattern  
572 Recognition Conference*, pages 29083–29095.

- 573 Yi Chen, Yuying Ge, Yixiao Ge, Mingyu Ding, Bohao  
574 Li, Rui Wang, Ruiyang Xu, Ying Shan, and Xihui  
575 Liu. 2023. Egoplan-bench: Benchmarking multi-  
576 modal large language models for human-level plan-  
577 ning. *arXiv preprint arXiv:2312.06722*. 577
- 578 Yilong Chen, Xiang Bai, Zhibin Wang, Chengyu Bai,  
579 Yuhan Dai, Ming Lu, and Shanghang Zhang. 2025b.  
580 Streamkv: Streaming video question-answering with  
581 segment-based kv cache retrieval and compression.  
582 *arXiv preprint arXiv:2511.07278*. 582
- 583 Alexandre Défossez, Laurent Mazaré, Manu Orsini,  
584 Amélie Royer, Patrick Pérez, Hervé Jégou, Edouard  
585 Grave, and Neil Zeghidour. 2024. Moshi: a speech-  
586 text foundation model for real-time dialogue. *arXiv  
587 preprint arXiv:2410.00037*. 587
- 588 Shangzhe Di, Zhelun Yu, Guanghao Zhang, Haoyuan  
589 Li, Tao Zhong, Hao Cheng, Bolin Li, Wanggui He,  
590 Fangxun Shu, and Hao Jiang. 2025. Streaming video  
591 question-answering with in-context video kv-cache  
592 retrieval. *arXiv preprint arXiv:2503.00540*. 592
- 593 Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch,  
594 Aakanksha Chowdhery, Ayzaan Wahid, Jonathan  
595 Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang,  
596 and 1 others. 2023. Palm-e: An embodied multi-  
597 modal language model. 597
- 598 Dave Epstein, Boyuan Chen, and Carl Vondrick. 2020.  
599 Oops! predicting unintentional action in video. In  
600 *Proceedings of the IEEE/CVF conference on com-  
601 puter vision and pattern recognition*, pages 919–929.  
602
- 602 Shenghao Fu, Qize Yang, Yuan-Ming Li, Yi-Xing Peng,  
603 Kun-Yu Lin, Xihan Wei, Jian-Fang Hu, Xiaohua Xie,  
604 and Wei-Shi Zheng. 2025. Vispeak: Visual instruc-  
605 tion feedback in streaming videos. *arXiv preprint  
606 arXiv:2503.12769*. 606
- 607 Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Neva-  
608 tia. 2017. Tall: Temporal activity localization via  
609 language query. In *Proceedings of the IEEE interna-  
610 tional conference on computer vision*, pages 5267–  
611 5275. 611
- 612 Eric Horvitz. 1999. Principles of mixed-initiative user  
613 interfaces. In *Proceedings of the SIGCHI conference  
614 on Human Factors in Computing Systems*, pages 159–  
615 166. 615
- 616 Zhenpeng Huang, Xinhao Li, Jiaqi Li, Jing Wang, Xi-  
617 angyu Zeng, Cheng Liang, Tao Wu, Xi Chen, Liang  
618 Li, and Limin Wang. 2025. Online video under-  
619 standing: Ovbench and videochat-online. In *Proceedings  
620 of the Computer Vision and Pattern Recognition Con-  
621 ference*, pages 3328–3338. 621
- 622 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam  
623 Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow,  
624 Akila Welihinda, Alan Hayes, Alec Radford, and 1  
625 others. 2024. Gpt-4o system card. *arXiv preprint  
626 arXiv:2410.21276*. 626

627	Jie Lei, Tamara L Berg, and Mohit Bansal. 2021. Detecting moments and highlights in videos via natural language queries. <i>Advances in Neural Information Processing Systems</i> , 34:11846–11858.	682
628		683
629		684
630		
631	Wei Li, Bing Hu, Rui Shao, Leyang Shen, and Liqiang Nie. 2025. Lion-fs: Fast & slow video-language thinker as online video assistant. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 3240–3251.	685
632		686
633		687
634		688
635		689
636	Junming Lin, Zheng Fang, Chi Chen, Zihao Wan, Fuwen Luo, Peng Li, Yang Liu, and Maosong Sun. 2024. Streamingbench: Assessing the gap for mllms to achieve streaming video understanding. <i>arXiv preprint arXiv:2411.03628</i> .	690
637		
638		
639		
640		
641	Zhenyu Ning, Guangda Liu, Qihao Jin, Wenchao Ding, Minyi Guo, and Jieru Zhao. 2025. Livevlm: Efficient online video understanding via streaming-oriented kv cache and retrieval. <i>arXiv preprint arXiv:2505.15269</i> .	691
642		692
643		693
644		694
645		695
646	Junbo Niu, Yifei Li, Ziyang Miao, Chunjiang Ge, Yuanhang Zhou, Qihao He, Xiaoyi Dong, Haodong Duan, Shuangrui Ding, Rui Qian, and 1 others. 2025. Ovobench: How far is your video-langs from real-world online video understanding? In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 18902–18913.	696
647		697
648		698
649		699
650		700
651		701
652		
653	Rui Qian, Shuangrui Ding, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Dahua Lin, and Jiaqi Wang. 2025. Dispider: Enabling video llms with active real-time interaction via disentangled perception, decision, and reaction. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 24045–24055.	702
654		703
655		704
656		705
657		706
658		707
659		
660	Rui Qian, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Shuangrui Ding, Dahua Lin, and Jiaqi Wang. 2024. Streaming long video understanding with large language models. <i>Advances in Neural Information Processing Systems</i> , 37:119336–119360.	708
661		709
662		710
663		711
664		
665	ByteDance Seed, Jiaze Chen, Tiantian Fan, Xin Liu, Lingjun Liu, Zhiqi Lin, Mingxuan Wang, Chengyi Wang, Xiangpeng Wei, Wenyuan Xu, and 1 others. 2025. Seed1. 5-thinking: Advancing superb reasoning models with reinforcement learning. <i>arXiv preprint arXiv:2504.13914</i> .	712
666		713
667		714
668		715
669		716
670		
671	Yemin Shi, Yu Shu, Siwei Dong, Guangyi Liu, Jaward Sesay, Jingwen Li, and Zhiting Hu. 2025. Voila: Voice-language foundation models for real-time autonomous interaction and voice role-play. <i>arXiv preprint arXiv:2505.02707</i> .	724
672		725
673		726
674		727
675		
676	Yansong Tang, Dajun Ding, Yongming Rao, Yu Zheng, Danyang Zhang, Lili Zhao, Jiwen Lu, and Jie Zhou. 2019. Coin: A large-scale dataset for comprehensive instructional video analysis. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 1207–1216.	728
677		729
678		
679		
680		
681		
682	Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. Bert rediscovers the classical nlp pipeline. <i>arXiv preprint arXiv:1905.05950</i> .	682
683		683
684		684
685		
686	Haibo Wang, Bo Feng, Zhengfeng Lai, Mingze Xu, Shiyu Li, Weifeng Ge, Afshin Dehghan, Meng Cao, and Ping Huang. 2025a. Streambridge: Turning your offline video large language model into a proactive streaming assistant. <i>arXiv preprint arXiv:2505.05467</i> .	685
687		686
688		687
689		688
690		689
691	Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan Chen, Yaohui Wang, and 1 others. 2023. Internvid: A large-scale video-text dataset for multimodal understanding and generation. <i>arXiv preprint arXiv:2307.06942</i> .	691
692		692
693		693
694		694
695		695
696	Yueqian Wang, Xiaojun Meng, Yuxuan Wang, Jianxin Liang, Jiansheng Wei, Huishuai Zhang, and Dongyan Zhao. 2024a. Videollm knows when to speak: Enhancing time-sensitive video comprehension with video-text duet interaction format. <i>arXiv preprint arXiv:2411.17991</i> .	696
697		697
698		698
699		699
700		700
701		701
702	Yuxuan Wang, Yueqian Wang, Bo Chen, Tong Wu, Dongyan Zhao, and Zilong Zheng. 2025b. Omnimmi: A comprehensive multi-modal interaction benchmark in streaming video contexts. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 18925–18935.	702
703		703
704		704
705		705
706		706
707		707
708	Yuxuan Wang, Cihang Xie, Yang Liu, and Zilong Zheng. 2024b. Videollamb: Long-context video understanding with recurrent memory bridges. <i>arXiv preprint arXiv:2409.01071</i> .	708
709		709
710		710
711		711
712	Cheng-Kuang Wu, Zhi Rui Tam, Chieh-Yen Lin, Yun-Nung Vivian Chen, and Hung-yi Lee. 2024a. Streambench: Towards benchmarking continuous improvement of language agents. <i>Advances in Neural Information Processing Systems</i> , 37:107039–107063.	712
713		713
714		714
715		715
716		716
717	Shiwei Wu, Joya Chen, Kevin Qinghong Lin, Qimeng Wang, Yan Gao, Qianli Xu, Tong Xu, Yao Hu, Enhong Chen, and Mike Zheng Shou. 2024b. Videollm-mod: Efficient video-language streaming with mixture-of-depths vision computation. <i>Advances in Neural Information Processing Systems</i> , 37:109922–109947.	717
718		718
719		719
720		720
721		721
722		722
723		723
724	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwén Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, and 1 others. 2025. The rise and potential of large language model based agents: A survey. <i>Science China Information Sciences</i> , 68(2):121101.	724
725		725
726		726
727		727
728		728
729		729
730	Haomiao Xiong, Zongxin Yang, Jiazu Yu, Yunzhi Zhuge, Lu Zhang, Jiawen Zhu, and Huchuan Lu. 2025. Streaming video understanding and multi-round interaction with memory-enhanced knowledge. <i>arXiv preprint arXiv:2501.13468</i> .	730
731		731
732		732
733		733
734		734
735	Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan, Kai Dang, and 1 others. 2025a. Qwen2. 5-omni technical report. <i>arXiv preprint arXiv:2503.20215</i> .	735
736		736
737		737
738		738

739	Jin Xu, Zhifang Guo, Hangrui Hu, Yunfei Chu, Xiong Wang, Jinzheng He, Yuxuan Wang, Xian Shi, Ting He, Xinfu Zhu, Yuanjun Lv, Yongqi Wang, Dake Guo, He Wang, Linhan Ma, Pei Zhang, Xinyu Zhang, Hongkun Hao, Zishan Guo, and 19 others. 2025b. Qwen3-omni technical report. <i>arXiv preprint arXiv:2509.17765</i> .	795
740		796
741		797
742		798
743		799
744		800
745		
746	Ruyi Xu, Guangxuan Xiao, Yukang Chen, Liuning He, Kelly Peng, Yao Lu, and Song Han. 2025c. Streamgvlm: Real-time understanding for infinite video streams. <i>arXiv preprint arXiv:2510.09608</i> .	801
747		802
748		803
749		804
750		
751	Haolin Yang, Feilong Tang, Lingxiao Zhao, Xiang An, Ming Hu, Huifa Li, Xinlin Zhuang, Yifan Lu, Xiaofeng Zhang, Abdalla Swikir, and 1 others. 2025a. Streamagent: Towards anticipatory agents for streaming video understanding. <i>arXiv preprint arXiv:2508.01875</i> .	805
752		806
753		807
754		808
755		
756	Pinci Yang, Xin Wang, Xuguang Duan, Hong Chen, Runze Hou, Cong Jin, and Wenwu Zhu. 2022. Avqa: A dataset for audio-visual question answering on videos. In <i>Proceedings of the 30th ACM international conference on multimedia</i> , pages 3480–3491.	809
757		810
758		811
759		812
760		813
761	Zhenyu Yang, Yuhang Hu, Zemin Du, Dizhan Xue, Shengsheng Qian, Jiahong Wu, Fan Yang, Weiming Dong, and Changsheng Xu. 2025b. Svbench: A benchmark with temporal multi-turn dialogues for streaming video understanding. <i>arXiv preprint arXiv:2502.10810</i> .	814
762		
763		
764		
765		
766		
767	Zhenyu Yang, Kairui Zhang, Yuhang Hu, Bing Wang, Shengsheng Qian, Bin Wen, Fan Yang, Tingting Gao, Weiming Dong, and Changsheng Xu. 2025c. Livestar: Live streaming assistant for real-world online video understanding. <i>arXiv preprint arXiv:2511.05299</i> .	815
768		816
769		817
770		818
771		819
772		820
773	Zhiwei Yang, Chen Gao, Jing Liu, Peng Wu, Guansong Pang, and Mike Zheng Shou. 2025d. Assistpda: An online video surveillance assistant for video anomaly prediction, detection, and analysis. <i>arXiv preprint arXiv:2503.21904</i> .	821
774		822
775		
776		
777		
778	Linli Yao, Yicheng Li, Yuancheng Wei, Lei Li, Shuhuai Ren, Yuanxin Liu, Kun Ouyang, Lean Wang, Shicheng Li, Sida Li, and 1 others. 2025. Timechat-online: 80% visual tokens are naturally redundant in streaming videos. In <i>Proceedings of the 33rd ACM International Conference on Multimedia</i> , pages 10807–10816.	823
779		824
780		825
781		826
782		
783		
784		
785	Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zihui He, and 1 others. 2024. Minicpm-v: A gpt-4v level mllm on your phone. <i>arXiv preprint arXiv:2408.01800</i> .	827
786		
787		
788		
789		
790	Haoji Zhang, Yiqin Wang, Yansong Tang, Yong Liu, Jiashi Feng, Jifeng Dai, and Xiaojie Jin. 2024a. Flash-vstream: Memory-based real-time understanding for long video streams. <i>arXiv preprint arXiv:2406.08085</i> .	828
791		
792		
793		
794		

## 827 A Appendix

### 828 A.1 Related Works

829 To address the fragmented landscape of streaming  
 830 multimodal models, we unify representative  
 831 methods in a comparative analysis along two axes:  
 832 supported input modalities and interaction capabili-  
 833 ties. We observe that many works described as  
 834 “streaming” in fact adopt a question-injection pro-  
 835 tocol, where a query is issued at a predetermined  
 836 timestamp and the model answers using only the  
 837 preceding context. As a result, they primarily study  
 838 long-horizon processing via KV-cache compres-  
 839 sion and external memory, rather than continuous  
 840 online interaction with response-timing decisions.  
 841 In contrast, the few systems that support online  
 842 streaming interaction typically span all three in-  
 843 teraction types. LiveCC is a notable exception: it  
 844 focuses on fixed-rate real-time narration and there-  
 845 fore does not require deciding when to respond.  
 846 Moreover, LION-FS, VideoLLM-MoD, and LiveS-  
 847 tar mainly introduce efficiency improvements on  
 848 top of the VideoLLM-online pipeline. Accord-  
 849 ingly, we use VideoLLM-online as the represen-  
 850 tative baseline. Overall, Table 10 shows that our  
 851 method is the first open-source model to enable full  
 852 omni-modal streaming while natively supporting  
 853 proactive response, real-time narration, and reactive  
 854 QA within a unified framework.

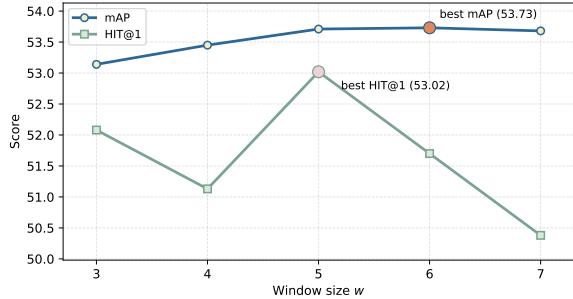
855 We also summarize commonly used benchmarks  
 856 for streaming evaluation in Table 1. Although these  
 857 benchmarks are often described as “streaming”,  
 858 they target different capabilities, and their coverage  
 859 is uneven, which motivates us to consolidate them  
 860 into a unified evaluation protocol.

### 861 A.2 Sensitivity Analysis

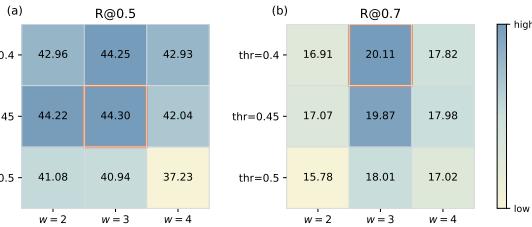
862 Sensitivity analysis confirms robust performance  
 863 (Figure 5). In static settings, mAP remains stable  
 864 while HIT@1 shows only slight sensitivity to varia-  
 865 tions in the window size. Dynamic tasks exhibit a  
 866 broad operating regime with smooth degradation,  
 867 indicating no brittle reliance on specific parame-  
 868 ters. Narration is likewise insensitive to speak head  
 869 probability thresholds, justifying a fixed default  
 870 without additional tuning (Table 8).

### 871 A.3 Evaluation Details

872 We specify evaluation protocols for our streaming  
 873 interaction tasks. For PO, we preprocess each sam-  
 874 ple to replicate the original benchmark by cropping  
 875 the video to the annotated ask time and injecting the



876 Figure 5: Sensitivity analysis on window size on  
 877 QVHighlight.



878 Figure 6: Sensitivity analysis on window size and thresh-  
 879 old on Charades-STA.

880 question at that timestamp, ensuring strictly causal  
 881 temporal ordering. While Streaming VLM base-  
 882 lines take text prompts, our model takes a speech  
 883 rendering of the same text to benchmark native  
 884 multimodal processing. For streaming baselines  
 885 (e.g., VideoLLM-Online and MMDuet), we record  
 886 the first-response timestamp and report accuracy as  
 887 the fraction of samples whose first-response time  
 888 is within  $\pm 2$  seconds of the annotated ground-truth  
 889 time. For REC, a model gets one point if its cho-  
 890 sen response time falls within the annotated event  
 891 interval. Each segment is evaluated once, and we  
 892 report the micro success rate. For CRR, we award  
 893 one point if the first response after the ask time  
 894 occurs after the annotated clue time, validating that  
 895 the model waits for necessary visual evidence be-  
 896 fore answering. Hyperparameters were determined  
 897 via validation sets as follows: QVHighlight win-  
 898 dow size = 5; Charades-STA window size = 3 with  
 899 threshold = 0.45; PA window size = 5 with thresh-  
 900 old = 0.5; PO window size = 4 with threshold = 0.2;  
 901 REC window size = 2 with threshold = 0.7; CRR  
 902 window size = 2 with threshold = 0.7; YouCook2  
 903 threshold = 0.975; SSR threshold = 0.97.

904 Due to space constraints, in the main table we  
 905 report only the average score on the narration task,  
 906 computed as the mean of the three GPT-4o-based  
 907 evaluation dimensions. We present the full break-  
 908 down in Table 9.

Probability	YouCook2			OVO-Bench (SSR)		
	Threshold	F1	BERTScore	GPT-Eval	F1	BERTScore
0.965	35.36	0.82	0.52 / 0.28 / 0.31	<b>15.53</b>	0.82	0.59 / 0.28 / 0.29
0.970	35.05	0.82	0.50 / 0.29 / 0.33	14.54	0.83	<b>0.59 / 0.33 / 0.34</b>
<b>0.975</b>	<b>35.21</b>	<b>0.83</b>	<b>0.53 / 0.29 / 0.36</b>	14.58	0.83	0.62 / 0.32 / 0.33
0.980	34.90	0.83	0.55 / 0.28 / 0.36	15.15	<b>0.84</b>	0.59 / 0.33 / 0.36
0.985	34.07	0.83	0.52 / 0.29 / 0.31	14.73	0.84	0.62 / 0.29 / 0.35

Table 8: Sensitivity analysis of the probability threshold on real-time narration, presenting performance metrics (F1, BERTScore, GPT-Eval) across different triggering thresholds on YouCook2 and OVO-Bench (SSR).

Method	YouCook2			SSR		
	F1	BERTScore	GPT-Eval	F1	BERTScore	GPT-Eval
TimeChat	21.70	–	–	–	–	–
VTG-LLM	17.50	–	–	–	–	–
VideoLLM-online	18.82	0.82	0.33 / 0.05 / 0.12	10.24	<b>0.84</b>	0.39 / 0.02 / 0.14
MMDuet	17.81	0.83	0.31 / 0.26 / 0.12	9.02	0.79	0.42 / 0.29 / 0.21
<b>Ours</b>	<b>35.21</b>	<b>0.83</b>	<b>0.53 / 0.29 / 0.36</b>	<b>14.54</b>	0.83	<b>0.59 / 0.33 / 0.34</b>
<i>- Ablation Study</i>						
Mixed Training	31.42	0.81	0.47 / 0.30 / 0.24	8.88	0.80	0.52 / 0.34 / 0.13
w/o Speak Head	9.25	0.79	0.32 / 0.32 / 0.09	3.39	0.77	0.41 / 0.30 / 0.08
$K = 1$	34.43	0.82	0.51 / 0.30 / 0.31	9.64	0.78	0.49 / 0.24 / 0.22
<i>- Sensitivity Analysis</i>						
$w_{pos} = 2$	27.82	0.83	0.62 / 0.27 / 0.46	10.38	0.57	0.63 / 0.21 / 0.29
$w_{pos} = 4$	35.55	0.81	0.64 / 0.27 / 0.49	13.48	0.75	0.54 / 0.20 / 0.28

Table 9: Streaming narration results on YouCook2 and OVO-Bench (SSR). We report F1 for temporal window alignment, and use BERTScore and GPT-4o scores to assess narration quality.

#### A.4 Data Construction Details

**Proactive Data Processing** Since the original temporal annotations in DiDeMo and Charades-STA are often coarse, simply using them for streaming supervision introduces noise. We therefore re-annotate event windows using Doubao-Seed-1.6-thinking (Seed et al., 2025) to obtain precise start and end timestamps. For timing supervision, we label every second within the refined ground-truth event window as a positive trigger, ensuring the model learns robust event sensitivity.

**Narration Data Processing** Raw videos often contain unlabeled gaps that induce hallucination during training. We mitigate this by excising these intervals and concatenating annotated segments into continuous, semantically dense streams with recalibrated timestamps. For timing supervision, we discard the broad-window labeling of prior work (e.g., MMDuetIT) in favor of strict transition-based triggering. This yields precise supervision for incremental narration via multi-turn SFT.

**Audio Query Synthesis** To simulate real-world interaction, we synthesize text queries via TTS and overlay them onto original audio tracks. We strictly align these spoken queries with streaming units to enforce audio-driven instruction following.

#### A.5 Implementation Details

**Training Configuration.** We sample videos at 2 FPS and resize frames to a maximum of 65,536 pixels. The model is trained using LLaMA-Factory (Zheng et al., 2024) with a sequence length of 32K on 32 H20 GPUs (using a global batch size of 512). Proactive samples are specifically formatted as multi-turn dialogues to handle multiple triggers within a single stream.

**Streaming Decoding Logic.** Following the pipelined setting, if a response exceeds the 25-token budget, we append an `<|endoftext|>` (`<eot>`) token to signal an unfinished utterance. Decoding resumes in the subsequent segment and terminates only when `<|im_end|>` is generated.

Method	Inputs	Alert	Narration	Reactive QA	Description
Moshi (Défossez et al., 2024)	A	✗	✗	✓	A full-duplex speech–text model enabling low-latency, real-time voice dialogue without vision.
Voila (Shi et al., 2025)	A	✗	✗	✓	Real-time voice-language model with expressive role-play and vocal styles.
VideoStreaming (Qian et al., 2024)	T+V	✗	✗	✓	Streaming framework with hierarchical memory for coherent long-video understanding.
Flash-VStream (Zhang et al., 2024a)	T+V	✗	✗	✓	Real-time system with lightweight memory for low-latency long video processing.
VideoLLaMB (Wang et al., 2024b)	T+V	✗	✗	✓	Long-context model with recurrent memory bridges to propagate information over time.
InternLM-XComposer2.5 (Zhang et al., 2024b)	T+V	✗	✗	✓	Multimodal streaming system with layered memory for long-term video–audio interaction.
StreamChat (Xiong et al., 2025)	T+V	✗	✗	✓	Streaming agent using hierarchical memory to sustain long-context, multi-round dialogue.
StreamBridge (Wang et al., 2025a)	T+V	✗	✗	✓	Plug-and-play buffer turning offline Video-LLMs into proactive streaming assistants.
CogStream (Zhao et al., 2025)	T+V	✗	✗	✓	QA framework keeping only memory-critical context for efficient reasoning.
ReKV (Di et al., 2025)	T+V	✗	✗	✓	Retrieves historical KV-cache as in-context memory instead of re-encoding past frames.
LiveVLM (Ning et al., 2025)	T+V	✗	✗	✓	Optimizes streaming KV-cache update/retrieval for efficient long-horizon processing.
StreamAgent (Yang et al., 2025a)	T+V	✗	✗	✓	Uses KV-based temporal memory to anticipate events and respond proactively.
StreamingVLM (Xu et al., 2025c)	T+V	✗	✗	✓	Manages rolling KV-cache to support infinite video streams under bounded computation.
StreamKV (Chen et al., 2025b)	T+V	✗	✗	✓	Segments and compresses KV-cache, keeping only salient past segments within budget.
LiveCC (Chen et al., 2025a)	T+V	✗	✓	✗	Trained with continuous speech transcription for real-time narration of long videos.
VideoLLM-online (Chen et al., 2024)	T+V	✓	✓	✓	Online model for temporally aligned dialogue via a LIVE training framework.
AssistPDA (Yang et al., 2025d)	T+V	✓	✓	✓	Surveillance assistant unifying anomaly prediction and interactive analysis.
LION-FS (Li et al., 2025)	T+V	✓	✓	✓	Uses fast–slow thinking with selective tokenization for efficient streaming.
VideoLLM-MoD (Wu et al., 2024b)	T+V	✓	✓	✓	Enables efficient streaming by letting each layer skip a subset of tokens directly.
LiveStar (Yang et al., 2025c)	T+V	✓	✓	✓	Uses perplexity-based timing and streaming-aware attention for proactive understanding.
MMDuet (Wang et al., 2024a)	T+V	✓	✓	✓	Adopts a video–text duet format to insert replies during continuous playback.
Dispider (Qian et al., 2025)	T+V	✓	✓	✓	Disentangles perception, decision, and reaction for asynchronous responses.
MiniCPM-o 2.6 (Yao et al., 2024)	T+V+A	✗	✗	✓	Omni model supporting low-latency streaming speech over text, audio, and video.
Qwen2.5-Omni (Xu et al., 2025a)	T+V+A	✗	✗	✓	Dense Thinker–Talker model with streaming encoders for real-time perception.
Qwen3-Omni (Xu et al., 2025b)	T+V+A	✗	✗	✓	Natively omni-modal MoE architecture for high-concurrency streaming generation.
Stream-Omni (Zhang et al., 2025)	T+V+A	✗	✗	✓	Efficient modality alignment enabling speech interaction grounded on visual inputs.
ViSpeak (Fu et al., 2025)	T+V+A	✓	✗	✓	Vision-centric framework producing instruction-like feedback from evolving streams.
<b>ROMA (Ours)</b>	T+V+A	✓	✓	✓	<b>End-to-end streaming OLLM processing time-aligned chunks for proactive interaction.</b>

Table 10: Comparison of different streaming multimodal methods. **Note:** T=Text, V=Visual, A=Audio. ✓ supports the ability, ✗ does not.

946      **A.6 Case Study**

947      **Event-Triggered Alert** We present two event-  
948      triggered alert cases: one where the target event  
949      occurs only once (Figure 7), and another where it  
950      recurs multiple times (Figure 8). Compared with  
951      several representative VideoLLMs, our model trig-  
952      gers at more accurate times.

953      **Narration** In the narration task, the model must  
954      choose when to speak during a long streaming  
955      video and provide concise summaries of events  
956      observed so far without access to future content.  
957      As shown in Figure 9, compared with VideoLLMs,  
958      our outputs are more succinct and our response tim-  
959      ings align more closely with key event boundaries,  
960      leading to more accurate online narration.

961      **Reactive QA** With audio queries in the reactive  
962      QA setting (Figure 10), ROMA correctly localizes  
963      the relevant segment in the long video and extracts  
964      the key visual evidence. In contrast, MiniCPM-o  
965      misidentifies the segment, while Qwen2.5-Omni  
966      often responds with unnecessary follow-up ques-  
967      tions.

968      **A.7 Evaluation Prompt**

969      LLM-as-a-judge is a widely adopted paradigm for  
970      scalable evaluation, given its strong alignment with  
971      human preferences (Zheng et al., 2023). Accord-  
972      ingly, we employ GPT-4o as a reliable scorer for  
973      our open-ended tasks. Detailed prompts are pro-  
974      vided in Figures 11 and 12.

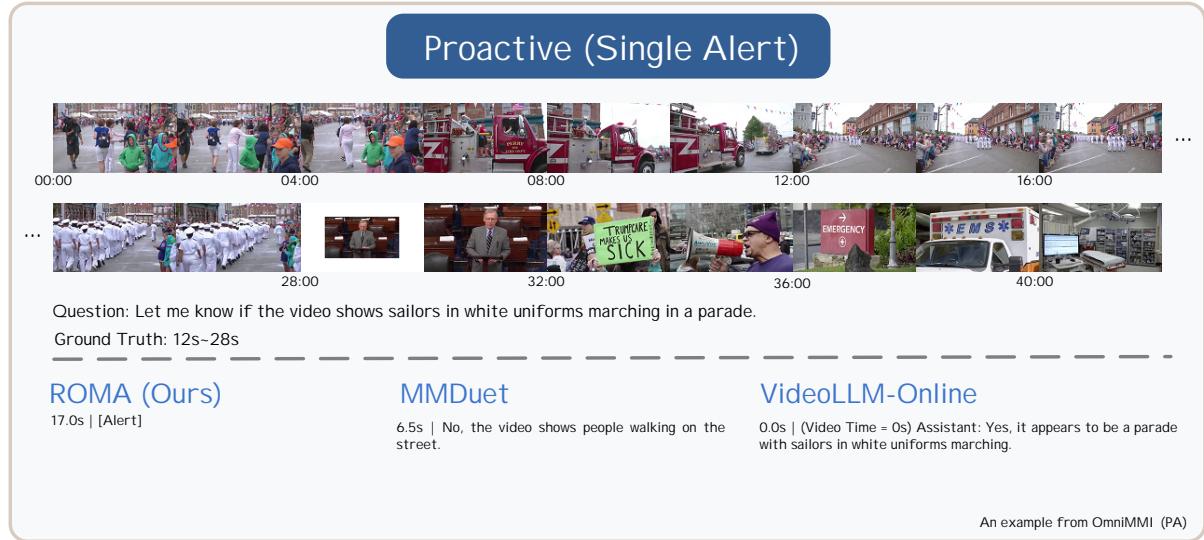


Figure 7: Qualitative comparison on the single-alert proactive task. While MMDuet and VideoLLM-Online exhibit premature triggering and hallucination before the target event appears, ROMA accurately accumulates visual evidence to release a precise alert at 17.0s, aligning with the ground truth interval (12s–28s).



Figure 8: Qualitative comparison on the recurring-alert task. While MMDuet suffers from continuous over-generation without distinguishing event boundaries, ROMA effectively tracks the repetitive action, releasing distinct alerts at 7.0s and 15.0s to capture the recurring instances.

## Narration



Question: Describe various activity events occurring in the video in real time.

Ground Truth:

- 6.0s | cut the goose liver into small slices
- 12.0s | season them with salt and pepper
- 22.0s | sprinkle with flour on both sides
- 30.0s | heat the apple in a pan with some oil
- 48.0s | sprinkle some salt on top of the apples and add some butter
- 81.0s | cook the liver on another pan with oil

### ROMA (Ours)

- 7.0s | Cut the liver into slices.
- 13.0s | Season the liver with salt and pepper.
- 24.0s | Sprinkle flour on both sides of the liver.
- 42.0s | Add apples to the pan and move them around.
- 59.0s | Add butter to the pan and shake the pan occasionally.
- 71.0s | Place the liver in the pan and cook until browned.

### MMDuet

- 5.0s | A person is cutting a piece of meat into slices on a cutting board.
- 11.0s | sprinkle salt and pepper
- 16.0s | sprinkle flour on the meat
- 21.0s | sprinkle flour on the meat
- 27.0s | add seasoning and flour to the meat
- 34.0s | pour oil into the pan
- 40.0s | fry meat
- 46.0s | fry the meat
- 52.0s | fry the meat
- 58.0s | add butter
- 64.0s | fry the meat
- 70.0s | add the ingredients to the pan
- 83.0s | put the liver into the pan
- 87.0s | fry meat

### VideoLLM-Online

- 0.0s | (Video Time = 0s) Assistant: You cut the meat with a knife.

An example from YouCook2

Figure 9: Qualitative comparison on the real-time narration task. While MMDuet suffers from severe repetition and redundant over-generation, ROMA effectively tracks the procedural evolution, generating concise, time-aligned descriptions that correspond strictly to the distinct ground truth events.

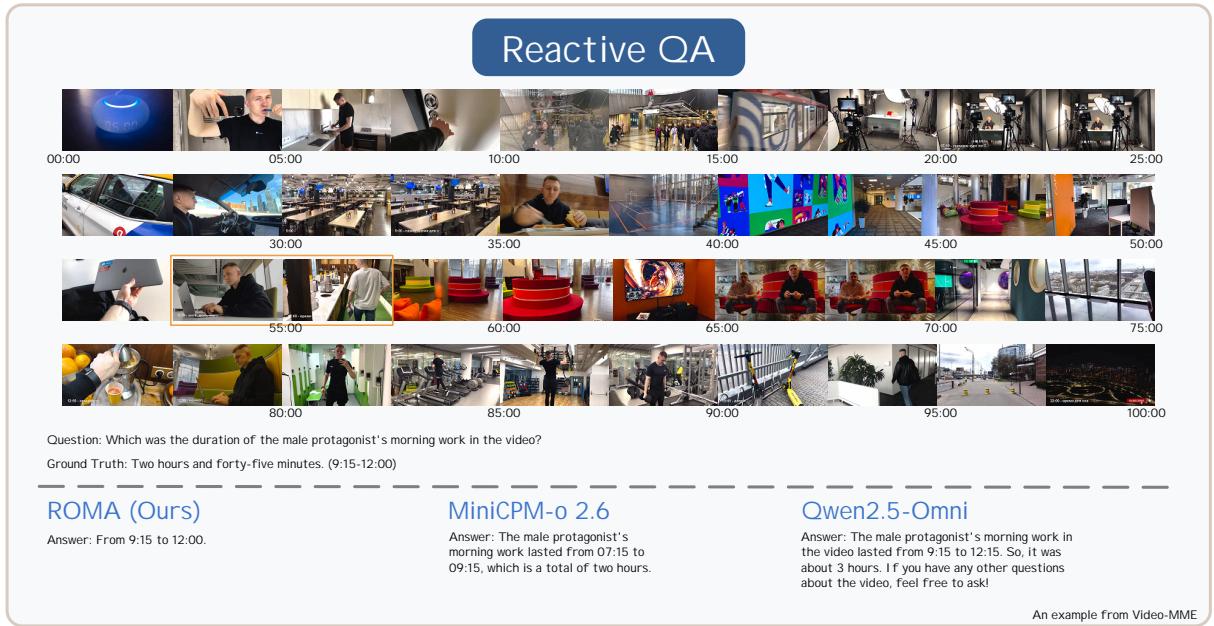


Figure 10: Qualitative comparison on the reactive QA task. While baseline models suffer from temporal misalignment or hallucinated intervals when querying specific activity durations, ROMA accurately retrieves the exact start and end timestamps (9:15–12:00) to derive the correct answer.

## Evaluation Prompt for Open-Ended Audio-Query Assessment

"""

You are an AI assistant tasked with evaluating whether a response matches the correct answer to a given question.

### ## Evaluation Rules

(1) Output 1 if the response matches the answer exactly or with synonymous/ equivalent wording.

- Synonyms, paraphrases, or different surface forms of the same meaning count as matches.

- Responses containing the correct answer with irrelevant details count as matches.

- Responses providing sufficient information to infer the correct answer count as matches.

(2) Output 0 if the response is incorrect, contradictory, or completely irrelevant to the question.

- If the answer and the response address different topics, or if the response does not answer the question.

- If the response introduces additional details that change the meaning of the answer, mark as 0.

### ### Examples

Example 1:

Question: What is the genre of this video?

Answer: It is a news report that introduces the history behind Christmas decorations.

Response: It's a Christmas-themed video, filled with festive decorations and a warm, cozy atmosphere. It really captures that classic holiday spirit. What do you think?

Your output: 0

Example 2:

Question: How many birds are above the fireplace?

Answer: 2.

Response: One bird is above the fireplace, and another is below.

Your output: 1

Your Turn:

Question: {question}

Answer: {ground\_truth}

Response: {prediction}

Your output:

"""

Figure 11: Full prompt provided to GPT-4o for open-ended evaluation with audio queries on Video-MME and EgoSchema.

## Evaluation Prompt for Narration Task Evaluation

"""

You are an expert evaluator for video narration quality. Your task is to compare a reference description of a video (ground truth) with a model-generated description for the same video, and output THREE scores between 0 and 1.

You must consider the model response as a SINGLE long story (it may contain multiple sentences describing different moments in the video).

IMPORTANT: Higher scores are always better.

Definitions:

1. coherence (story coherence):

- How internally coherent and well-structured is the model-generated story by itself?
- Does it read like a reasonable, temporally plausible sequence of actions and states?
- Penalize contradictions, abrupt jumps, and incoherent, rambling structure.
- 1.0 = very coherent and well-structured; 0.0 = completely incoherent.

2. alignment (semantic alignment with ground truth):

- How well does the model-generated story capture the key actions and steps in the ground truth?
- Consider whether important actions/events are present, correctly described, and roughly in a reasonable order.
- Hallucinated major steps that clearly do not appear in the ground truth should reduce this score.
- 1.0 = almost all key content in GT is covered with correct semantics; 0.0 = almost completely unrelated.

3. conciseness (relevant non-redundancy / brevity):

- This score measures whether the model response is concise GIVEN IT IS RELEVANT to the ground truth.
- If the model response is largely unrelated to the ground truth (low semantic overlap, wrong topic, ignores the video), conciseness MUST be near 0, even if the response is short.
- Penalize heavy repetition of similar sentences, long irrelevant digressions, and obvious padding.
- However, do NOT penalize necessary detail that genuinely helps describe the steps.
- 1.0 = succinct, minimal redundancy while preserving essential details;
- 0.0 = extremely repetitive / rambling / full of irrelevant filler / irrelevant with the groundtruth.

Empty or meaningless model responses (or responses that ignore the task) should receive low scores,  
typically near 0 for all dimensions.

Output format (VERY IMPORTANT):

- You MUST output valid JSON with exactly the following keys:  
  {"coherence": <float>, "alignment": <float>, "concreteness": <float>}
- Each value must be a number between 0 and 1 (inclusive).
- Do NOT output any extra text or explanation.

"""

Figure 12: Prompt used to instruct GPT-4o to evaluate narration quality along three criteria: story coherence, alignment with ground truth, and conciseness.