

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

Shehreen Azad¹, Vibhav Vineet², Yogesh S Rawat¹

¹ Center for Research in Computer Vision, University of Central Florida;

² Microsoft Research

CVPR 2026

Streaming Video Understanding



Causal understanding

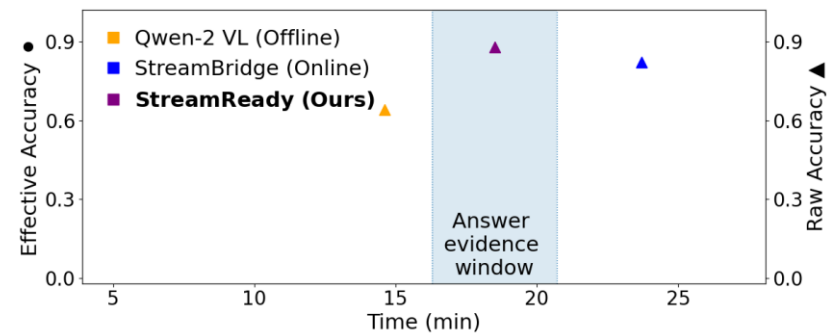
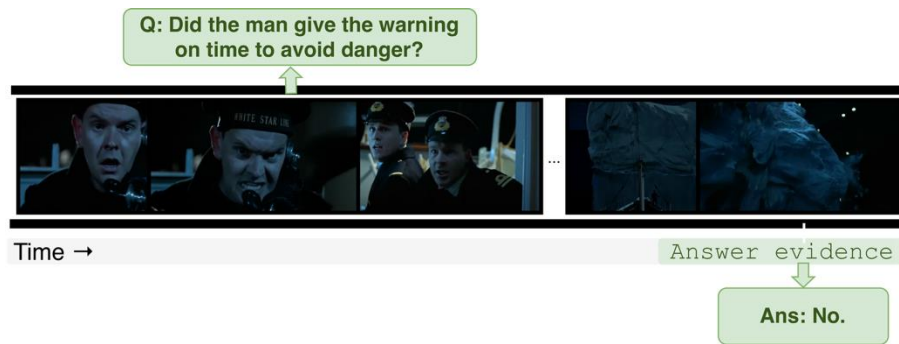
- Current or past-dependent
- Example: “Which team’s player started the match?”

Proactive understanding

- Future-dependent
- Example: “Output GOAL when a goal has been scored.”
- When is there enough evidence to answer?

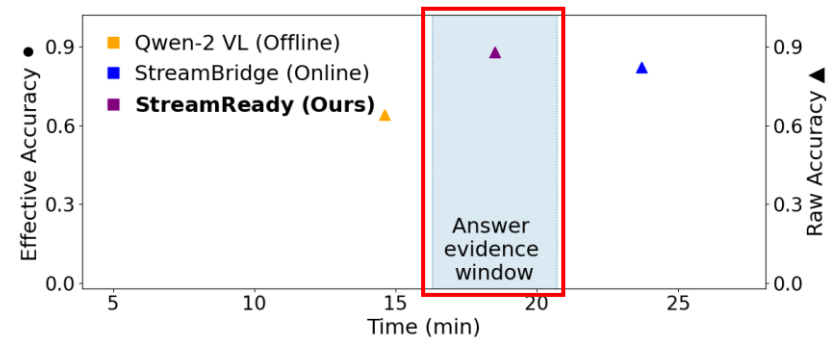
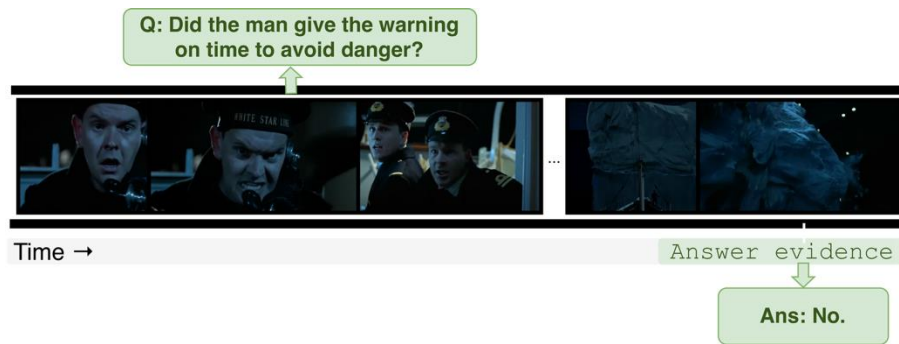
Current Methods Ignore Answer Timing

- Answer content focused → ignore timing.
- No notion of evidence sufficiency.
- Cannot distinguish between mistimed and on-timed answer.

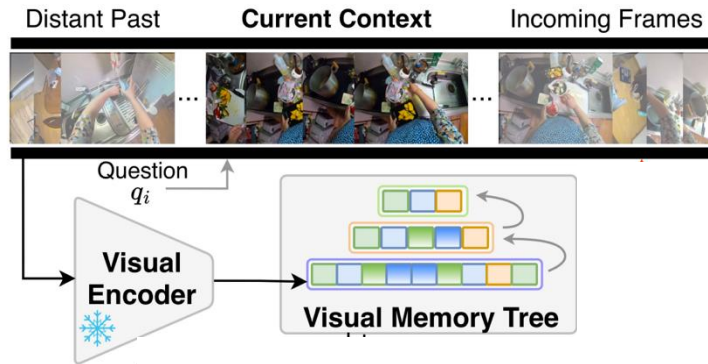


Towards Timing-Aware Answering

- Readiness-aware streaming understanding.
- Models should answer on-time based on visual evidence.



Task-Agnostic Storage

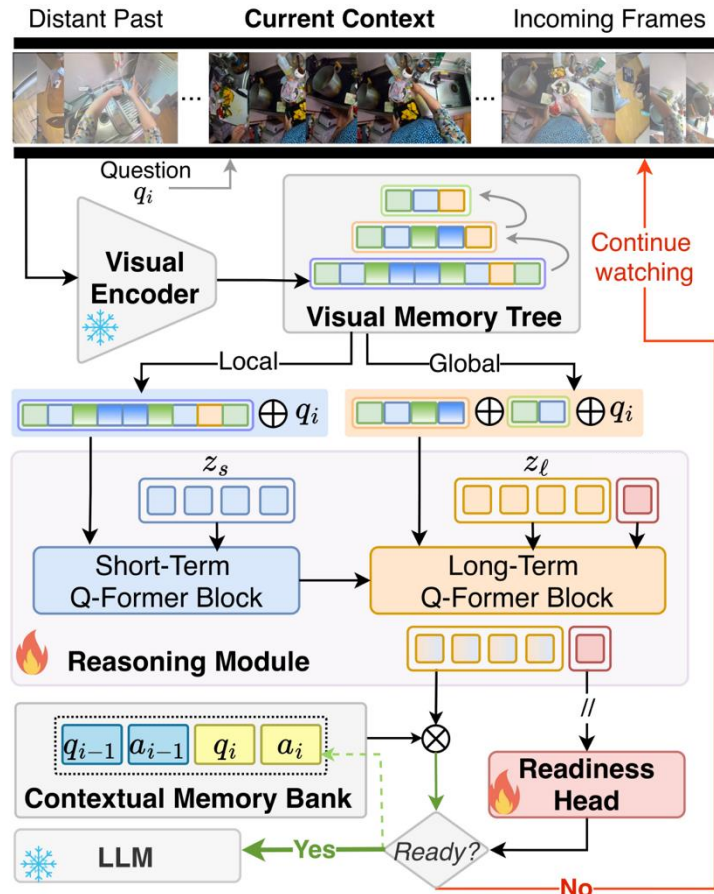


Visual Memory Tree

- Multi-level with increasing abstraction.
- Lowest level: short term details.
 - FIFO update.
 - Evicted frames \rightarrow next level.
- Higher levels: long term summary
 - Initialized by clustering (K-means)
 - Updated with EMA
- When question appears \rightarrow task-aware reasoning.

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

Readiness Monitoring



- Introduce learnable token in long-term reasoning block of HierarQ.
- Readiness head monitors token and assigns high score when sufficient evidence.
- If score exceeds threshold, trigger LLM.

Learning Readiness

- No answer evidence ground truth
- Similarity to memory and current answer representation.
 - High similarity \rightarrow pseudo-positive, Low similarity \rightarrow pseudo-negative.
- Optimize with contrastive loss.



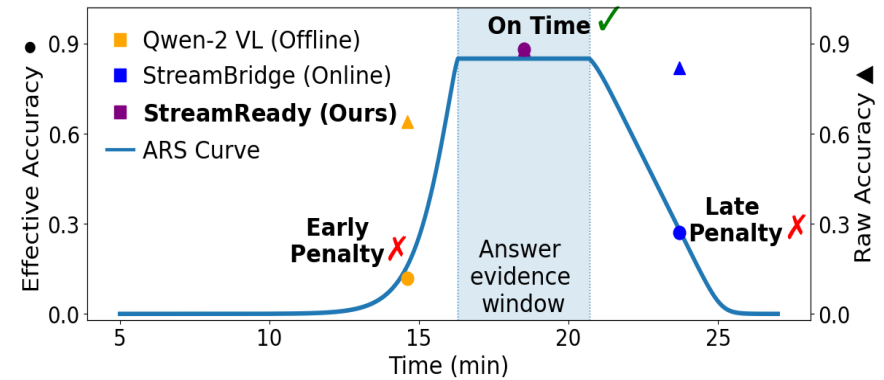
$$\mathcal{L}_{ctr} = -\log \sigma (R_{pred}(t^+) - R_{pred}(t^-))$$

Answer Readiness Score (ARS)

- Timing-aware evaluation metric to reward on-time answers

$$ARS = \frac{1}{N} \sum_{i=1}^N (EP_i \cdot LP_i); \quad \in [0, 1]$$

- Harsh early penalty
- Mild late penalty
- Defines effective accuracy



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

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

$$Acc_e = Acc \times ARS$$

Proactive Readiness Benchmark

Sequential Step Recognition	Repetitive Event Counting
 <p>Q. Time: 0:0; Q: Describe the steps of making vegetable stir-fry. Ans. Time: 0; Ans: < step 1> → Time: 1:21; Ans: < step 2 > → Time: 5:20; Ans: < step 3 >..</p>	 <p>Q Time: 1:10; Q: How many trees did the man pick up? Time: 1:24 (1) → Time: 1:44 (2) → Time: 1:57 (3, 4). Ans. Time: 1:57 Ans: 4</p>
Clues Reveal Responding	Causal Trigger Detection
 <p>Q Time: 15:10; Q: The woman is preparing to cook. Tell me when she starts putting things in the pot? Ans. Time: 35:05; (The woman puts spices in pot) Ans: Now.</p>	 <p>Q Time: 5:10; Q: The man is taking measurements. When is it known if he did it correctly? Ans. Time: 10:00; (The ply is shorter than the space.) Ans: 10:00</p>
Goal-State Detection	
 <p>Q Time: 12:20; Q: The man hands the girl a drink. Tell me when she finishes drinking. Ans: Time: 15:34; Ans: Finished</p>	

ProReady-QA

- 5 proactive understanding tasks with answer evidence window.
- 30 min – 1 hour long videos.
- 5000 QA pairs.

Results

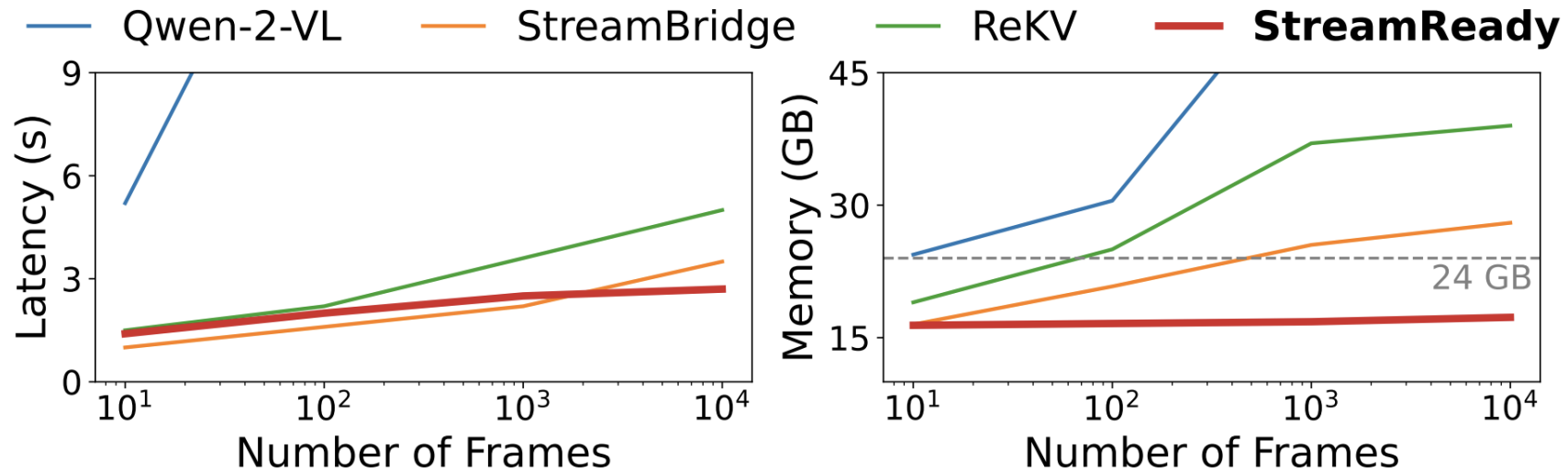
- 9 benchmarks spanning ~2 min to ~1 hour.
- Task: VQA (streaming understanding, offline understanding)
- Evaluation metric: Accuracy, ARS, Effective accuracy

Method	StreamingBench (1-10 min)	ProReady-QA (30 min – 1 h)		
	Acc	Acc	ARS	Eff. Acc
HierarQ <small>(CVPR'25)</small>	35.1	46.0	0.40	27.0 Δ -19.0
Previous Best <small>StreamBridge(NeurlPS'25)</small>	43.9	53.1	0.60	42.3 Δ -10.8
StreamReady	48.2 <small>+4.3</small>	56.4 <small>+3.3</small>	0.69 <small>+0.09</small>	53.2 Δ -3.2

Contribution of Each Component

Method	ProReady-QA	
	Acc.	ARS
Baseline	28.7	0.44
+ Task-Agnostic Storage	32.4 _{+3.7}	0.46 _{+0.02}
+ Task-Aware Reasoning	39.4 _{+7.0}	0.48 _{+0.02}
+ Readiness Monitoring	39.6 _{+0.2}	0.68 _{+0.20}

Cost Analysis



Lightweight readiness monitoring system and task-agnostic storage keeps cost minimum.

Thank You