ReWind: Understanding Long Videos with Instructed Learnable Memory

Anxhelo Diko^{1*†} Tinghuai Wang^{2*} Wassim Swaileh² Shiyan Sun² Ioannis Patras²

¹La Sapienza University of Roma ²Huawei Helsinki Research Center

diko@di.uniromal.it {tinghuaiwang, shiyansun, wassim.swaileh, ioannis.patras}@huawei.com

setting

VQA:

dataset: MovieChat-1K dataset (1000 for 9.13min)

1. global: requires processing the entire video and answering questions about its content

2. breakpoint: processing the video up to a specific timestamp and answering questions about the event at that point

metric: acc; score

Temporal Grounding:

dataset: Charades-STA

metric: mIoU; R@XX

motivation

Challenges:

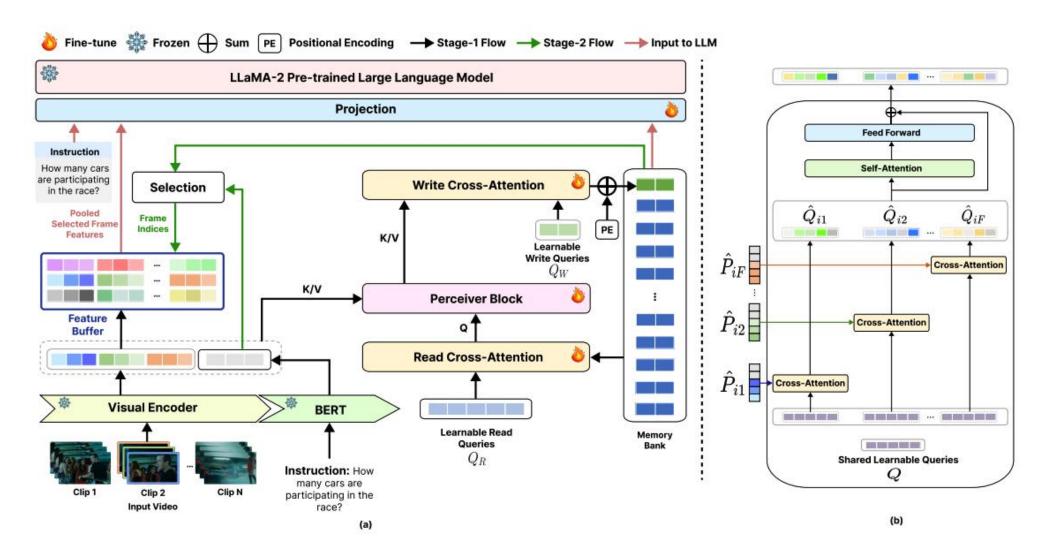
1. self-attention mechanisms require substantial memory that scales quadratically with the number of tokens, making long video processing computationally intensive

use Q-Former structure to compress spatial information

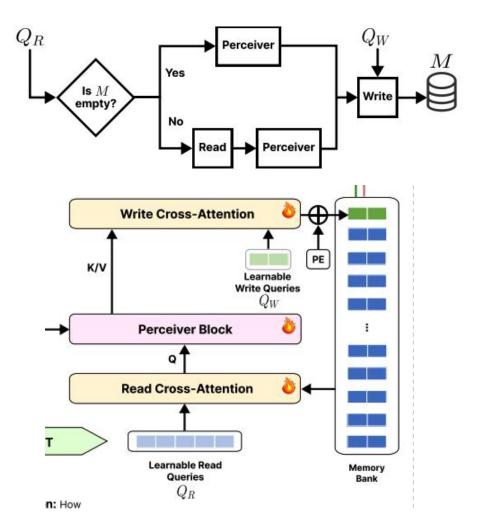
2. struggle to effectively model temporal dependencies over extended sequences

use memory-informed method to preserve temporal fidelity

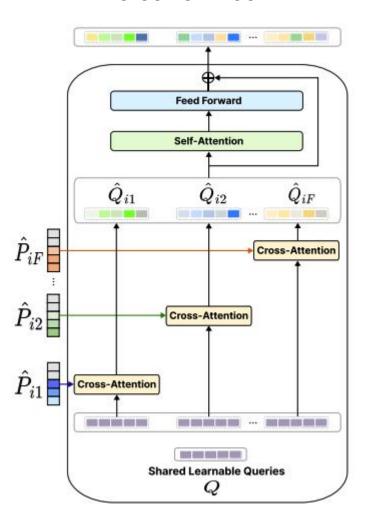
overview



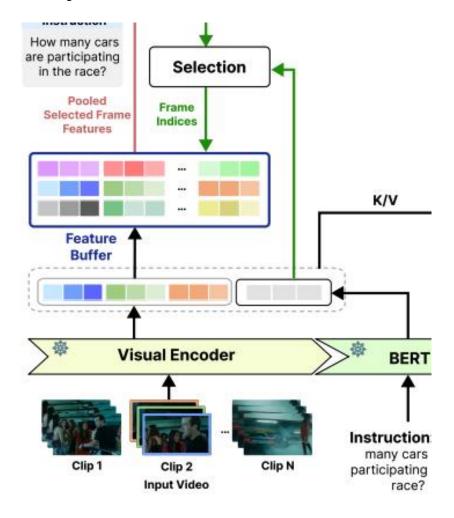
Instructed Memory Architecture



Perceiver Block



Dynamic Frame Selection



 Instruction based selection compute the attention matrix between
 I and contents of M select top L frames set Z_I
 Clustering

$$\sigma_l = exp(-\frac{1}{K} \sum_{z_k \in KNN(z_l, Z)} ||z_k - z_l||^2)$$

$$ho_l = egin{cases} \min_{j:\sigma_j > \sigma_l} ||z_j - z_l||^2 & \textit{if } \exists j \textit{ s.t. } \sigma_j > \sigma_l, \ \max_j \ ||z_j - z_l||^2 & \textit{otherwise}. \end{cases}$$

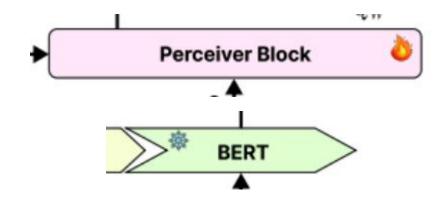
select top K (σ I \times ρ I) frames features Z

final input: <M, spec, Z>

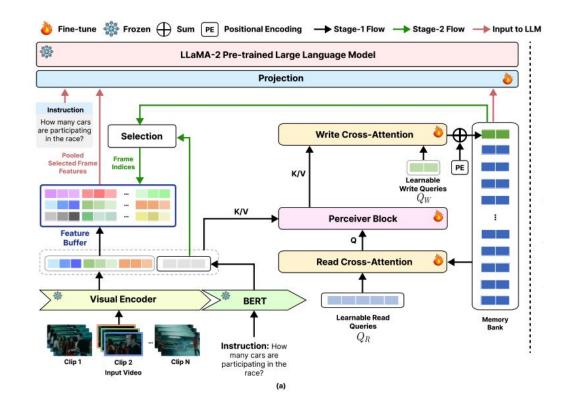
Training

Multimodal Pretraining Stage: contrastive learn

$$\mathcal{L}_{SigLIP} = -rac{1}{||\mathcal{B}||} \sum_{i=1}^{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} rac{\log(rac{1}{1 + \exp(z_{ij}(-t\mathbf{x}_i \cdot \mathbf{y}_j + b))})}{\mathcal{L}_{ij}})$$



Instruction Tuning Stage:



experiment

long video QA

Model	Num Frames	Num Tokens	Global VQA Break		Breakpoir	nt VQA Glob		Globa	obal Generation		Breakpoint Generation			n		
1710461	rum Frumes	Truin Tokens	Accuracy	Score	Accuracy	Score	CI	DO	CU	TU	co	CI	DO 2 2.85 1 3.32 5 3.09 7 3.24	CU	TU	co
Video LLaMA [27]	32	32	51.4	3.10	38.2	2.31	3.30	2.53	3.28	2.77	3.42	2.42	2.85	2.87	2.00	2.87
Video-ChatGPT [18]	100	356	44.2	2.71	49.8	2.71	2.48	2.78	3.03	2.48	2.99	3.11	3.32	3.29	2.62	3.29
Video Chat [14]	32	3072	61.0	3.34	48.3	2.43	3.26	3.20	3.38	2.97	3.47	2.96	3.09	3.24	2.46	3.22
MovieChat [21]	2048	8192	67.8	3.81	50.4	2.96	3.32	3.28	3.44	3.06	3.48	3.07	3.24	3.31	2.70	3.45
ReWind (Ours)	548*	1184*	80.6	4.46	57.2	3.4	4.18	4.00	4.24	4.02	3.54	3.41	3.37	3.64	2.97	3.61

short video QA

Model	LLM	Backbone	CI	DO	CU	TU	CO	AVG
Video Chat [14]	Vicuna-7B	ViT-G	2.23	2.50	2.53	1.94	2.24	2.29
Video LLaMA [27]	Vicuna-7B	ViT-G	1.96	2.18	2.16	1.82	1.79	1.98
Video-ChatGPT [18]	Vicuna-7B	ViT-L	2.40	2.52	2.62	1.98	2.37	2.38
LLaMA Adapter [28]	LLaMA-7B	ViT-L	2.03	2.32	2.30	1.98	2.15	2.16
Chat-UniVi [12]	Vicuna1.5-7B	ViT-L	2.89	2.91	3.46	2.39	2.81	2.89
VTimeLLM [11]	Vicuna1.5-7B	ViT-L	2.78	3.10	3.40	2.49	2.47	2.85
MovieChat [21]	LLaMA2-7B	ViT-G	2.76	2.93	3.01	2.24	2.42	2.67
LLaMA-VID [15]	Vicuna-7B	ViT-G	2.96	3.00	3.53	2.46	2.51	2.89
ReWind (Ours)	LLaMA2-7B	ViT-G	2.91	2.85	3.42	2.71	2.68	2.91

Temporal grounding

Model	Charades-STA							
Model	R@0.3	R@0.5	R@0.7	mIoU				
Video Chat [14]	9.0	3.3	1.3	6.5				
Video LLaMA [27]	10.4	3.8	0.9	7.1				
Video-ChatGPT [18]	20.0	7.7	1.7	13.7				
GroundingGPT [16]	-	29.6	11.9	2				
TimeChat [20]	-	32.2	13.4	-				
VTimeLLM [11]	51.0	27.5	11.4	31.2				
ReWind (Ours)	59.0	41.6	20.53	39.3				

AdaCM²: On Understanding Extremely Long-Term Video with <u>Adaptive</u> <u>Cross-Modality Memory Reduction</u>

Yuanbin Man¹, Ying Huang¹, Chengming Zhang², Bingzhe Li³, Wei Niu⁴, Miao Yin^{1†}

¹Department of CSE, University of Texas at Arlington, ²Department of CS, University of Houston, ³Department of CS, University of Texas at Dallas, ⁴School of Computing, University of Georgia

{yuanbin.man, ying.huang}@uta.edu, czhang48@uh.edu, bingzhe.li@utdallas.edu, wniu@uga.edu, miao.yin@uta.edu

setting

Long-Trem Video Understanding

dataset: LVU

content understanding and metadata prediction tasks

Video Captioning:

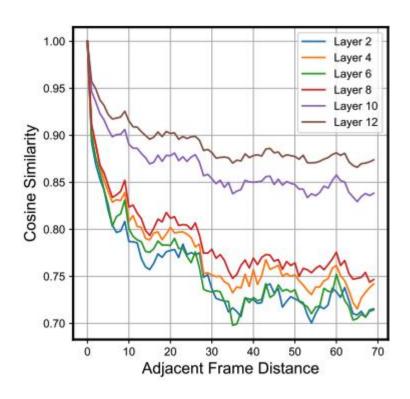
dataset: MSRVTT; MSVD; YouCook2

metric: METEOR"语义正确性" / CIDEr"表述符合共识"

motivation

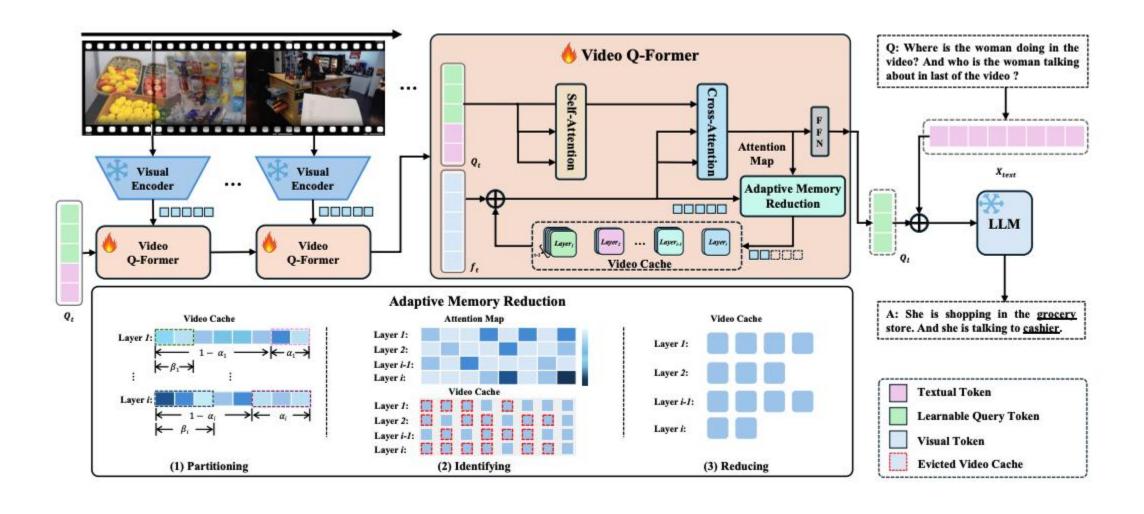
Observation

- 1. Only a subset of visual tokens exhibits high correlations to the text query within a frame
- 2. Correlation varies across different layers

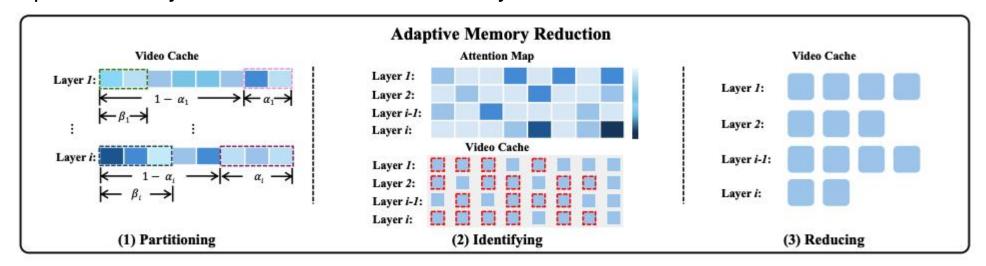


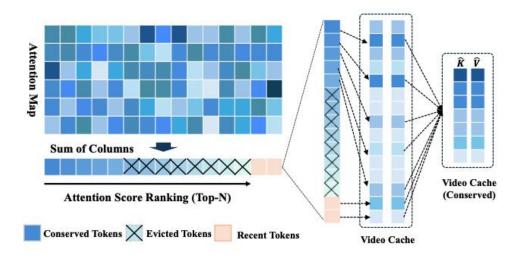
Adaptive Memory Reduction with Cross-Modality Attention

overview



Adaptive Memory Reduction with CrossModality Attention





Algorithm 1 Layer-wise Adaptive Video Reduction.

- 1: **Input:** Video frames $\{f_t\}_{t=1}^T$, hyper-parameters α, β for each layer;
- 2: Output: Reduced video KV cache K_t, V_t .
- 3: for t=1 to T do
- 4: $K_t \leftarrow [K_{t-1}, f_t W_K]; V_t \leftarrow [V_{t-1}, f_t W_V];$
- 5: **Partitioning:** $K_t \rightarrow [\hat{K}_t, \tilde{K}_t]; V_t \rightarrow [\hat{V}_t, \tilde{V}_t];$
- 6: **Identifying:** Determine S_t^c using Eq.5;
- 7: **Reducing:** Obtain \overline{K}_t using Eq.7;
- 8: $K_t \leftarrow [\overline{K}_t, \tilde{K}_t]; V_t \leftarrow [\overline{V}_t, \tilde{V}_t];$
- 9: end for

experiment

Table 1. Comparison with state-of-the-art methods on the LVU dataset. The <u>underlined</u> number means the second best.

especialists are asked as		Content			Metadata					
Model	Relation	Speak	Scene	Director	Genre	Writer	Year	Avg		
Obj_T4mer [45]	54.8	33.2	52.9	47.7	52.7	36.3	37.8	45.1		
Performer [9]	50.0	38.8	60.5	58.9	49.5	48.2	41.3	49.6		
Orthoformer [32]	50.0	38.3	66.3	55.1	55.8	47.0	43.4	50.8		
VideoBERT [36]	52.8	37.9	54.9	47.3	51.9	38.5	36.1	45.6		
LST [19]	52.5	37.3	62.8	56.1	52.7	42.3	39.2	49.0		
VIS4mer [19]	57.1	40.8	67.4	62.6	54.7	48.8	44.8	53.7		
S5 [43]	67.1	42.1	73.5	67.3	65.4	51.3	48.0	59.2		
MA-LMM [16]	58.2	44.8	80.3	74.6	61.0	<u>70.4</u>	<u>51.9</u>	63.0		
Ours	63.1	40.2	86.2	75.4	68.0	77.0	62.5	67.5		

Model	Breakfast	COIN
TSN [44]	121	73.4
VideoGraph [18]	69.5	(=)
Timeception [17]	71.3	-
GHRM [13]	75.5) = 0
D-Sprv. [25]	89.9	90.0
ViS4mer [19]	88.2	88.4
S5 [43]	90.7	90.8
MA-LMM [16]	<u>93.0</u>	<u>93.2</u>
Ours	94.4	93.3

ARREST 10000 CAST	MSRVTT		MS	SVD	YouCook2		
Model	M	C	M	C	M	C	
UniVL [28]	28.2	49.9	29.3	52.8	-	127.0	
SwinBERT [24]	29.9	53.8	41.3	120.6	15.6	109.0	
GIT [41]	32.9	73.9	51.1	180.2	17.3	129.8	
mPLUG-2 [47]	34.9	80.3	48.4	165.8	-	-	
VideoCoca [49]	-	73.2	-	-	-	128.0	
VideoLLaMA [52]	32.9	71.6	49.8	175.3	16.5	123.7	
MA-LMM [16]	33.4	<u>74.6</u>	51.0	179.1	17.6	131.2	
Ours	33.0	73.1	51.4	189.4	17.6	125.6	