

VidCtx: Context-aware Video Question Answering with Image Models

¹ Information Technologies Institute (ITI), Centre for Research and Technology Hellas (CERTH), Thessaloniki, Greece ² School of Electronic Engineering and Computer Science, Queen Mary University of London, London, UK {agoulas, bmezaris}@iti.gr, i.patras@qmul.ac.uk

Zero-shot Video Question Answering

Motivation

- Video LLMs achieve state-of-the-art Video QA performance
- High demand for **memory** and **compute**
- O(N²) space and time complexity
- Typically require large-scale video pre-training

Proposed method: VidCtx



- Integrate both text and visual modalities
- Process video frame-by-frame with image models
- Concatenate helpful context with prompts and visual features
- Aggregate decisions across frames with a pooling layer

Experimental results

Experimental setup

- Datasets: NExT-QA (5K), IntentQA (2K), STA
- **Base Model:** LLaVa-1.6-Mistral-7B

Comparison of VidCtx with zero-shot openapproaches

- **Competitive performance** on all datasets
- Scales linearly w.r.t. the number of frames
- 7B model even outperforms some models rely on proprietary LLMs (e.g. LLoVi with G





Andreas Goulas^{1,2}, Vasileios Mezaris¹, Ioannis Patras²



Captioner Prompt: Please provide a short description of the image giving information related to the following question: What did the white dog do after he looked up?

QA Prompt: Here is what happens earlier (or later) in the video: <Caption>. Question: What did the white dog do after he looked up? Option A: hit cans. Option B: yellow toy. Option C: walking. Option D: get up. Option E: smells the black dog. Option F: No Answer. Considering the information presented in the caption and the video frame, select the correct answer in one letter from the options (A,B,C,D,E,F).

Architecture Question-aware Captions: pre-trained image LLN generates captions related to the given question • Context-aware QA: image LLM answers question based on a) the <u>current frame</u> and b) the <u>caption</u> <u>a distant frame</u> (N/2 frames apart) **Aggregation:** frame-level decision scores are

	-		NExT-QA				STAR					
	Model	Params	Cau.	Tem.	Des.	All	IntentQA	Int.	Seq.	Pre.	Fea.	Avg
	GPT-based Models											
	LLoVi (GPT-3.5) [9]	N/A	67.1	60.1	76.5	66.3	-	-	-	-	-	-
AR (7K)	LLoVi (GPT-4) [9]	N/A	73.7	70.2	81.9	73.8	67.1	-	-	-	-	-
	VideoTree (GPT-4) [7]	N/A	75.2	67.0	81.3	73.5	66.9	-	-	-	-	-
	VideoAgent (GPT-4) [8]	N/A	72.7	64.5	81.1	71.3	-	-	-	-	-	-
	VFC [27] *	<1B	51.6	45.4	64.1	51.5	-	-	-	-	-	-
	InternVideo [28] *	< 1B	48.0	43.4	65.1	49.1	-	43.8	43.2	42.3	37.4	41.6
-model	SeViLA [6] *	4B	61.5	61.3	75.6	63.6	60.9	48.3	45.0	44.4	40.8	44.6
	LangRepo [10]	12B	64.4	51.4	69.1	60.9	59.1		-	-	-	-
	Q-ViD [11]	12B	<u>67.6</u>	<u>61.6</u>	72.2	<u>66.3</u>	63.6	48.2	47.2	43.9	43.4	45.7
	VideoChat2 [1] *	7B	61.9	57.4	69.9	61.7	-	62.4	67.2	57.5	53.9	63.8
	VidCtx (Ours)	7B	71.7	65.1	79.2	70.7	67.1	<u>53.9</u>	<u>54.3</u>	<u>51.4</u>	<u>44.7</u>	<u>51.1</u>
	* VFC [27] pre-train their model on 0.5M video-caption pairs. SeViLA [6] trains a localizer component on 10K											
_	videos. InternVideo [28]] and Vide	oChat2	[1] utiliz	ze large-	scale pr	e-training o	n video	datasets	with ov	ver 10M	
S	video-caption pairs.											
						C	omparison	with ca	ptions	-only ba	aseline (N	NExT-O
s that	Effect of number of frames (NExT-QA)											
	# Frames 1	2 4	8	16	32	64	$\frac{1}{Can}$	tions Onl	v (32 fr	ames)	67.3	
ירו אנ ² .5-ו אנ	Top-1 (%) 63.5 66	6.8 68.8	69.2	70.1	70.3	70.7	Vid	Ctx (32 f	rames)	11105)	70.3	
,								X				•



Frame 52 of 64

Caption: The white dog (...) appears to have continued lying on the floor .



Goals

aggregated with max pooling and L1 normalization

Ablation study

- Normalization prior to pooling is important

Choice of context (NExT-	QA)	Choice of aggregation method (NExT-Q/		
Method (use of context)	Top-1 (%)	Method (aggregation)	Top-1 (%)	
No Context	67.9	Voting	69.7	
Concat 16 Captions (Q-Aware)	68.3	Mean Pooling	69.3	
Current Caption (Q-Aware)	69.9	Max Pooling	69.2	
Distant Caption (Static)	69.5	Softmax + Mean Pooling	70.2	
Distant Caption (Q-Aware)	70.7	Softmax + Max Pooling	70.6	
		L1 Norm + Max Pooling	70.7	



he looked up? (Category: Temporal)



A) hit cans B) yelow toy C) walking D) get upE) smells the black dog F) No Answer

No C	ontext	Vid	Ctx
F20	F52	F20	F52
0.14	0.14	0.11	0.12
0.23	0.22	0.14	0.15
0.18	0.19	0.11	0.13
0.22	0.20	0.53	0.32
0.18	0.20	0.03	0.16
0.02	0.03	-0.06	-0.09

Address scaling and memory limitations of Video LLMs • Adopt a training-free paradigm by re-using pre-trained image LLMs and <u>combine visual-text representations</u>

Λ	$y = \underset{t \in T - \{F\}}{\operatorname{argmax}} \left[\underset{0 \le i < N}{\max} \frac{p(d_i = t)}{\sum_{k \in T} p(d_i = k) } \right]$
<u>of</u>	 Frame-level scores Extract scores from logits of first token
)n	 Ignore special No Answer token

• **Distant captions** provide the most helpful context • +3% improvement over captions-only baseline • Performance **improves** with higher **number of frames**



Supported by: