

The Devil is in Temporal Token: High Quality Video Reasoning Segmentation

Supplementary Material

001 This supplementary material provides additional details
 002 and analysis of VRS-HQ, expanding on the content pre-
 003 sented in the main paper. We begin by evaluating the im-
 004 pact of various training datasets on segmentation perfor-
 005 mance (§A). Next, we present more detailed implementa-
 006 tion information to facilitate reproducibility (§B). We then
 007 elaborate on the specific method of utilizing SAM2 [8]
 008 for mask decoding and propagation (§C). Subsequently, we
 009 show some failure cases with analysis to offer a more com-
 010 prehensive understanding of VRS-HQ’s limitations (§D).
 011 Then we compare VRS-HQ with other methods on the mul-
 012 timodal question-answering tasks. (§E) Additionally, we
 013 present more qualitative comparisons against VISA, high-
 014 lighting the strengths of our proposed method (§F). Finally,
 015 we visualize the reasoning segmentation results of VRS-HQ
 016 on in-the-wild video datasets, demonstrating its strong gen-
 017 eralization capabilities (§G).

018 A. Datasets Ablation

019 As illustrated in Tab. 1, fine-tuning with the full datasets
 020 yields the best performance while excluding the image
 021 segmentation dataset, VideoQA dataset [6], or ReVOS
 022 dataset [9] individually results in varying degrees of metric
 023 degradation. Notably, removing the VideoQA dataset mini-
 024 mally impacts the model’s performance, with a decline of
 025 **0.9%** in $\mathcal{J}\&\mathcal{F}$ on both the referring and reasoning subsets,
 026 as its primary role is to support the MLLM’s video com-
 027 prehension rather than directly contributing to the segmen-
 028 tation process. In contrast, excluding the ReVOS dataset
 029 leads to a noticeable drop of **4.4%** and **7.6%** in $\mathcal{J}\&\mathcal{F}$, high-
 030 lighting its pivotal role in enhancing the model’s reasoning
 segmentation performance in challenging scenarios.

Table 1. Ablation study on the impact of training datasets.

Datasets	referring			reasoning		
	\mathcal{J}	\mathcal{F}	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}	$\mathcal{J}\&\mathcal{F}$
Joint	59.8	64.5	62.1	53.5	58.7	56.1
w/o ImageSeg	58.5	63.2	60.8	51.0	56.3	53.6
w/o VideoQA	<u>58.7</u>	<u>63.7</u>	<u>61.2</u>	<u>52.4</u>	<u>58.0</u>	<u>55.2</u>
w/o ReVOS	55.3	60.1	57.7	45.3	51.6	48.5

031

032 B. Additional Implementation Details

033 Due to space constraints of the main document, additional
 034 implementation details are provided here. During training,
 035 we use varying sampling ratios for different datasets (*cf.*
 036 Tab. 2). For video segmentation datasets, 8-12 frames are
 037 uniformly sampled at fixed intervals per video, and up to

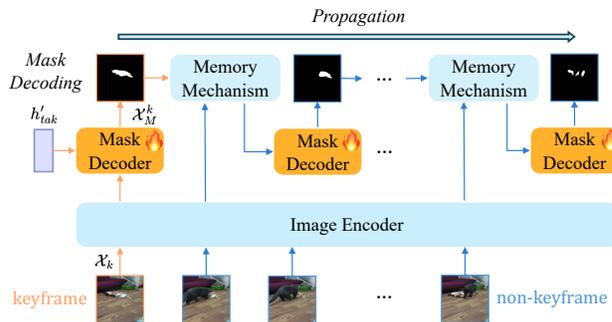


Figure 1. Details of SAM2 for mask decoding and propagation. All the video frames are input into the image encoder for feature extraction. The feature embeddings of the keyframe interact with h'_{tak} through the mask decoder for mask generation and then propagate it to the remaining video frames via the memory mechanism.

three object categories are selected per image or video. Dur- 038
 ing inference, we utilize CLIP-336 [7] for global sampling, 039
 selecting up to 12 frames per video. Input images are re- 040
 sized to 224×224 before being input to Chat-UniVi [3]. 041
 Data passed to SAM2 is augmented as described in [4] and 042
 resized to 1024×1024 . Moreover, LoRA [2] is applied 043
 with a scaling factor of 16 and a dropout rate of 0.05 across 044
 all query and value projection layers within the MLLM, en- 045
 abling efficient fine-tuning.

Table 2. Datasets sampling ratio during training.

Dataset	SemSeg	RIS	ImageQA	ReaSeg	VideoQA	VideoSeg
Ratios	9/32	3/32	3/32	1/32	1/8	3/8

046

047 C. More Details of SAM2

As depicted in Fig. 1, we provide detailed insights into the 048
 process of mask decoding and propagation using SAM2 [8]. 049
 Specifically, all input video frames are processed through 050
 the image encoder to extract multi-scale visual features. 051
 Subsequently, the fused temporal embedding h'_{tak} interacts 052
 with the keyframe features in the mask decoder to gener- 053
 ate the segmentation mask and perform video-level propa- 054
 gation. The prediction is then encoded by the memory en- 055
 coder and stored in the memory bank, which maintains a 056
 FIFO queue of memories from recent frames. Feature em- 057
 beddings from subsequent non-keyframes attend to these 058
 stored mask features through memory attention and utili- 059
 ze the mask decoder to generate corresponding masks, en- 060
 abling inter-frame propagation. 061

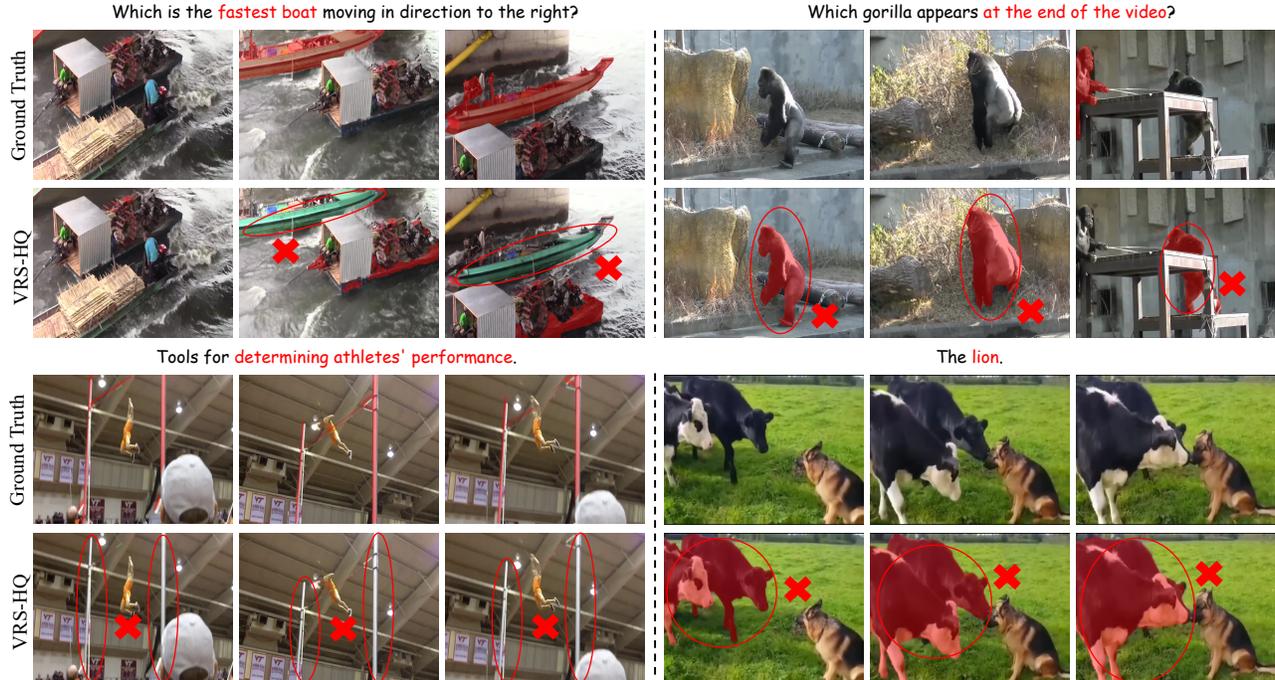


Figure 2. Visualization of failure cases for VRS-HQ. These examples illustrate the model’s limitations in scenarios requiring complex world knowledge and temporal reasoning, as well as challenges in processing negative samples.

062 D. Failure Case Analysis

063 Fig. 2 presents a detailed analysis of several failure cases, offering a deeper understanding of the limitations of VRS-HQ. 064 The **top row** highlights two specific challenges. **First**, 065 VRS-HQ struggles with keyframe localization when presented with queries based on motion, such as identifying the 066 fastest-moving boat within a video sequence. This suggests a potential weakness in analyzing and interpreting dynamic 067 visual information. **Second**, the model exhibits difficulty 068 segmenting targets with minimal temporal presence, as exemplified by the gorilla visible only in the last two frames 069 of the video. This points to a possible limitation in effectively capturing and utilizing short-duration visual cues. The **bot-** 070 **tom row** reveals further limitations. VRS-HQ demonstrates 071 a lack of comprehension when faced with nuanced or implicitly phrased prompts, such as recognizing a “high bar” 072 within the context of gymnastics performance evaluation. This suggests a need for improved understanding of complex 073 semantic relationships within video content. Furthermore, the model occasionally exhibits hallucinatory behavior, 074 generating segmentations for non-existent objects, particularly when dealing with empty targets or scenes where 075 the requested object is absent. 076 077 078 079 080 081 082 083 084

085 We hypothesize that several strategies could mitigate 086 these limitations. Improving the video comprehension 087 capabilities of the Multimodal Large Language Model

(MLLM) could enhance the ability to interpret complex 088 scenes and queries. Enabling the model to process a larger 089 number of sampled frames simultaneously might improve 090 its sensitivity to subtle temporal changes and short-duration 091 events. Finally, designing specialized tokens specifically for 092 representing empty masks could address the observed hal- 093 lucinations in such scenarios. We leave a thorough investi- 094 gation of these potential improvements to future research. 095

096 E. VQA Task Results

097 To explore the relationship between dense prediction tasks 098 and multimodal QA, we evaluate VRS-HQ against its founda- 099 tion MLLM Chat-UniVi, on the POPE [5] benchmark, as 100 illustrated in Tab. 3. VRS-HQ performs better on multi- 101 modal QA despite being designed for reasoning segmen- 102 tation, demonstrating the synergistic relationship between 103 these tasks and the potential for cross-task improvements.

Table 3. Results on the POPE benchmark.

Methods	POPE-R		POPE-P		POPE-A	
	Accuracy	F1-Score	Accuracy	F1-Score	Accuracy	F1-Score
Chat-UniVi	85.19	86.05	69.50	74.39	64.97	71.54
VRS-HQ	87.25	87.18	75.40	77.38	70.40	73.97

104 F. More Qualitative Comparison

105 In addition to the visual comparisons presented in the main 106 document, we provide further comparisons across more di-

107 verse settings in Fig. 3-5 to demonstrate the model’s reason-
108 ing and segmentation capabilities. As illustrated in Fig. 3,
109 VISA demonstrates reduced sensitivity to color-related ex-
110 pressions (e.g., “white” and “brown”) when provided with
111 explicit textual instructions. Furthermore, the example on
112 the left demonstrates VISA’s tendency to misidentify visu-
113 ally similar objects with complex spatial variations. In con-
114 trast, VRS-HQ effectively aggregates temporal information,
115 capturing inter-frame motion dynamics and leading to im-
116 proved segmentation accuracy.

117 Fig. 4 highlights the robust segmentation and reasoning
118 capabilities of VRS-HQ in scenarios with complex tempo-
119 ral dynamics. In the left example, VISA struggles to pre-
120 cisely detect the airplane appearing on the left at the end
121 of the video. Similarly, in the right case, VISA misclas-
122 sifies the tiger emerging in the lower left corner. In con-
123 trast, VRS-HQ leverages the Token-driven Keyframe Se-
124 lection for more accurate keyframe identification and inte-
125 grates SAM2 with the temporal token, enriched with both
126 intra-frame spatial and inter-frame temporal relations, re-
127 sulting in reliable decoding and consistent object tracking.

128 Fig. 5 presents scenarios requiring general and world
129 knowledge for reasoning. In the first example (left), VISA
130 segments only two koi carp (*Cyprinus carpio*) correctly,
131 whereas VRS-HQ identifies nearly all the fish present. In
132 the second example (right), VISA fails to associate “dog”
133 with the phrase “common household pet”, indicating limi-
134 tations in its reasoning capabilities. By contrast, VRS-HQ
135 leverages the integration of temporal tokens to achieve a
136 more nuanced semantic understanding, enabling finer con-
137 trol and interpretation.

138 G. In-the-wild Visualization Results

139 Fig. 6 and Fig. 7 show qualitative results of VRS-HQ on in-
140 the-wild videos. Fig. 6 shows results on first-person videos
141 from the GTEA dataset [1], using implicit prompts. Even
142 in cluttered kitchen environments with many similar ob-
143 jects, VRS-HQ demonstrates strong generalization capabil-
144 ity. It is particularly effective at segmenting smaller tar-
145 gets, such as the spoon and watch shown in the first and
146 third rows, respectively, maintaining robust performance in
147 these challenging scenarios. Fig. 7 shows results on 360-
148 degree panoramic videos from the PanoVOS dataset [10],
149 using more intricate prompts. Notably, VRS-HQ success-
150 fully segments individuals even when they are split across
151 the distorted edges of the video (first row), without any task-
152 specific optimizations. Furthermore, it maintains effective
153 tracking performance when the primary subjects within the
154 video are moving dynamically (last two rows).

References

- [1] Alireza Fathi, Xiaofeng Ren, and James M Rehg. Learning to recognize objects in egocentric activities. In *CVPR 2011*, pages 3281–3288. IEEE, 2011.
- [2] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [3] Peng Jin, Ryuichi Takanobu, Wancai Zhang, Xiaochun Cao, and Li Yuan. Chat-univi: Unified visual representation empowers large language models with image and video understanding. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 13700–13710, 2024.
- [4] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Int. Conf. Comput. Vis.*, pages 4015–4026, 2023.
- [5] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023.
- [6] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*, 2023.
- [7] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. pages 8748–8763. PMLR, 2021.
- [8] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.
- [9] Cilin Yan, Haochen Wang, Shilin Yan, Xiaolong Jiang, Yao Hu, Guoliang Kang, Weidi Xie, and Efstratios Gavves. Visa: Reasoning video object segmentation via large language models. In *Eur. Conf. Comput. Vis.*, 2024.
- [10] Shilin Yan, Xiaohao Xu, Renrui Zhang, Lingyi Hong, Wen-chao Chen, Wenqiang Zhang, and Wei Zhang. Panovos: Bridging non-panoramic and panoramic views with transformer for video segmentation. In *European Conference on Computer Vision*, pages 346–365. Springer, 2025.

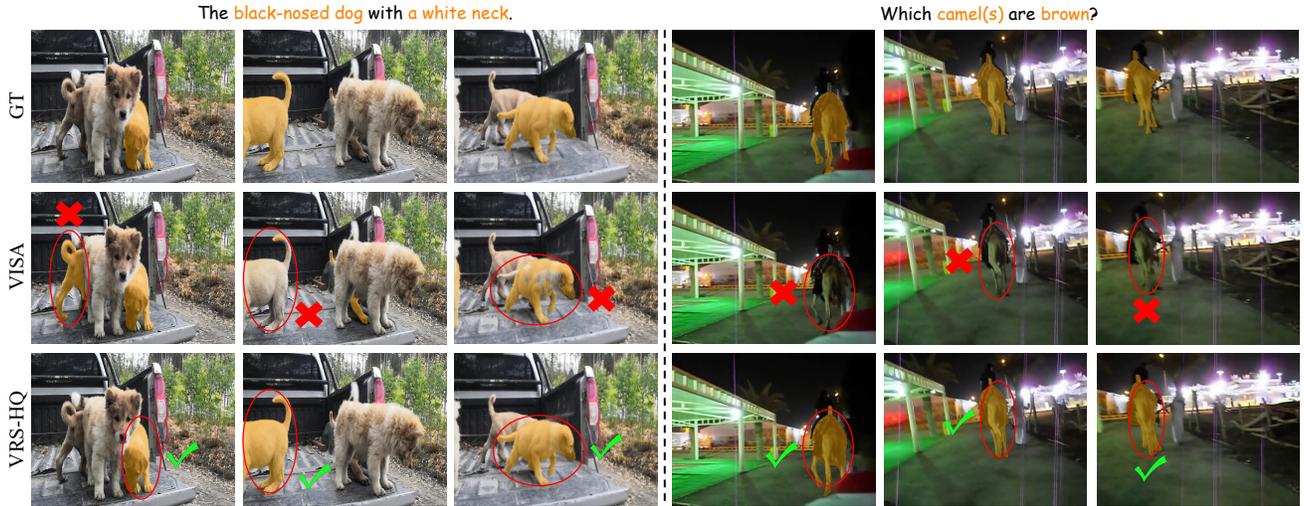


Figure 3. Qualitative comparison of VRS-HQ and VISA in explicit language-based referring scenarios on the ReVOS benchmark.

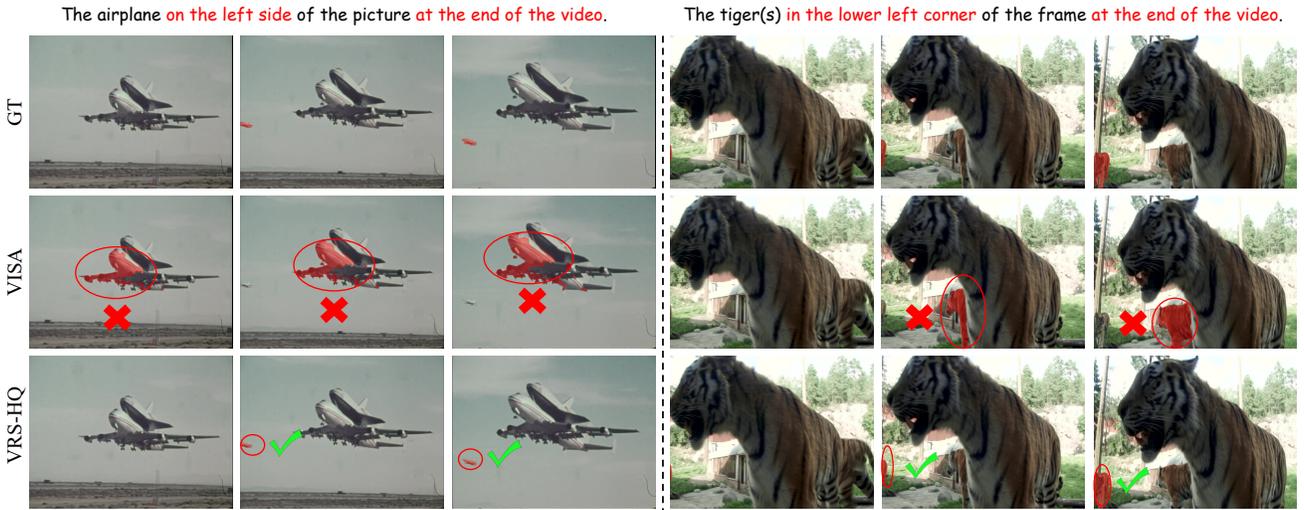


Figure 4. Qualitative comparison of VRS-HQ and VISA in scenarios incorporating complex temporal dynamics on the ReVOS benchmark.

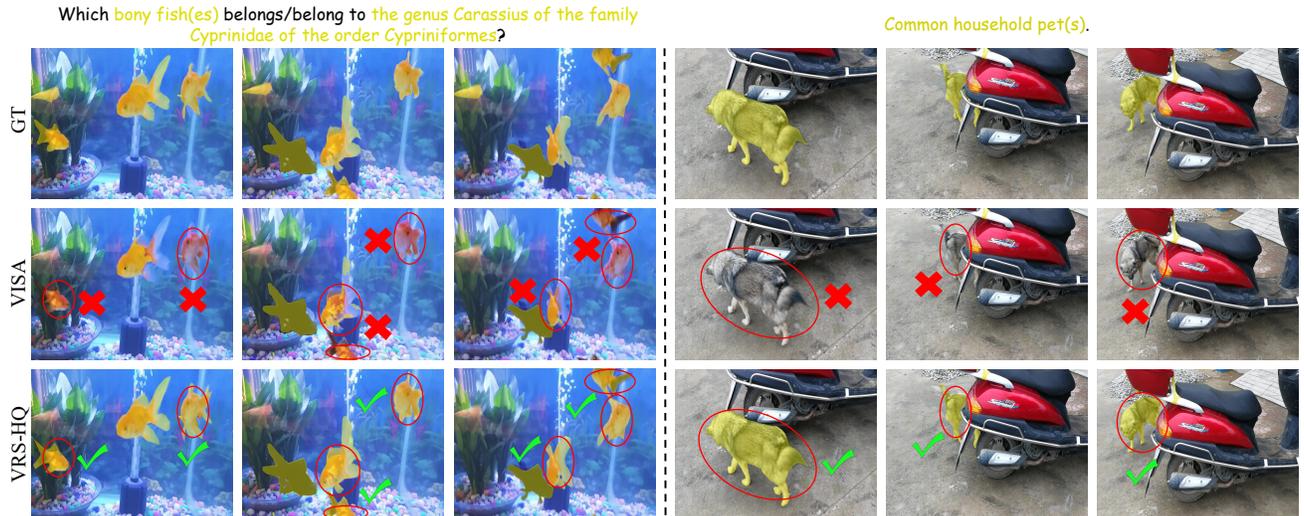


Figure 5. Qualitative comparison of VRS-HQ and VISA in reasoning scenarios that require world knowledge on the ReVOS benchmark.

The video's subject uses a curved tool for scooping and serving liquids or soft solids, which typically features a rounded bowl and a handle, ensuring ease of use and practicality.



A white object with a circular smooth surface, typically placed on a table and unmoved by the videographer, designed to hold and present various types of food, offering functionality in dining settings.



An item used for timekeeping in daily life and offering aesthetic appeal, commonly worn on the wrist and often featuring a circular or rectangular design, occasionally appearing in the video.



Figure 6. Visualization of VRS-HQ utilized in egocentric videos.

A prominent figure subtly highlighted within a commercial promotional video, whose presence and actions serve as the central point for engagement and communication.



The individual demonstrating a more seasoned presence and refined mastery of BMX Freestyle maneuvers, embodying the skill and experience often associated with years of dedication to the sport.



Who displays the most dynamic and expressive range of movement during the dance, transitioning seamlessly between sharp, high-energy motions, captivating attention with their vibrant and energetic performance.



Figure 7. Visualization of VRS-HQ applied to 360-degree panoramic videos.