

SAM3-DMS: Decoupled Memory Selection for Multi-target Video Segmentation of SAM3

Ruiqi Shen¹ Chang Liu^{2,✉} Henghui Ding^{1,✉}
¹Fudan University ²Shanghai University of Finance and Economics

<https://github.com/FudanCVL/SAM3-DMS>



Figure 1. SAM3 vs. SAM3-DMS (Ours), in simultaneous multi-object video segmentation. In Frame #3, SAM3 updates the memories for all objects together based on the overall status, saving Object 1 into memory even its mask is blank, causing identity drifts. In contrast, we separately update memory for each object by its own, keeping consistent identity tracking.

Abstract

Segment Anything 3 (SAM3) has established a powerful foundation that robustly detects, segments, and tracks specified targets in videos. However, in its original implementation, its group-level collective memory selection is suboptimal for complex multi-object scenarios, as it employs a synchronized decision across all concurrent targets conditioned on their average performance, often overlooking individual reliability. To this end, we propose SAM3-DMS, a training-free decoupled strategy that utilizes fine-grained memory selection on individual objects. Experiments demonstrate that our approach achieves robust identity preservation and tracking stability. Notably, our advantage becomes more pronounced with increased target density, establishing a solid foundation for simultaneous multi-target video segmentation in the wild.

1. Introduction

The Segment Anything Model (SAM) family has revolutionized pixel-wise visual perception, with SAM1 [14] and SAM2 [22] introducing Promptable Visual Segmentation (PVS) for images and videos, while SAM3 [2] advances Promptable Concept Segmentation (PCS) by detecting and tracking instances of high-level semantic concepts, marking a significant step toward high-level video grounding.

Built upon SAM2, SAM3 employs a memory bank to preserve target context across frames. To handle complex scenarios such as temporary disappearance and re-entry, SAM3 implements a “memory selection” strategy which conditionally updates the memory bank by thresholding the prediction confidence, memorizing only the reliable features. However, in simultaneous multi-target video segmentation scenarios, this selection strategy remains limited to coarse frame-level decision. Specifically, by aggregating confidence scores of distinct targets into a group-level average, SAM3 enforces a synchronized update decision

✉ Corresponding to: hhding@fudan.edu.cn, liuchang@sufe.edu.cn

across all concurrent targets driven by their collective performance, rather than the reliability of each individual.

While this approach suffices for single-object video segmentation, it may become problematic in real-world videos where the scene contains multiple targets, each with distinct and complex patterns, such as PCS cases. Fig. 1 shows a demonstration where Object #1 disappears from the scene in Frame #3. However, SAM3 overlooks this absence as the high confidence of other objects (*e.g.*, #2) elevates the group-level average score. Consequently, a blank mask is encoded into and polluted Object #1’s memory, causing identity drift at its subsequent re-entry.

To maintain temporal identity consistency of individual targets while adhering to the simultaneous paradigm of tracking and segmenting all targets at once, we propose SAM3-DMS, a simple yet effective training-free approach. Firstly, to address the group-level memory selection issue, we employ an decoupled policy where the selection is decided on each object independently, rather than the average of all concurrent targets. Secondly, the decision is subject to the target’s self-assessment. In each frame, each target derives a confidence score by combining its segmentation score with the overall visibility status, which serves as a criterion to determine whether or not to update its own corresponding memory. By eliminating distractions from other unrelated targets in the group, this fine-grained memory selection strategy ensures robust multi-target video segmentation in challenging scenarios involving occlusion, distractors, and target re-entries, as illustrated in Fig. 1.

Furthermore, we perform comprehensive evaluations for our decoupled memory selection strategy on both PCS and PVS tasks across seven benchmarks and in-the-wild videos under the simultaneous multi-target video segmentation setting, which demonstrate the superiority of our approach, notably revealing that the performance gap widens as target density increases. In summary, this work enhances SAM3’s multi-target robustness driven by fine-grained memory control. By decoupling memory maintenance conditioned on individual reliability, SAM3-DMS effectively mitigates identity drift and secures temporal consistency in challenging open-world scenarios.

2. Related Work

Memory-based VOS. Memory-based networks have become the dominant paradigm in Video Object Segmentation (VOS) [20], with the Space-Time Memory (STM) [19] serving as the foundational framework. Recent advances include STCN [6] for optimized attention computation, XMem [5] for decoupling long- and short-term memory, RMNet [26] with optical flow assistance, and Cutie [7] for object-level guidance. SAM2 [22] and SAM3 [2] have emerged as the state-of-the-art VOS methods. Trained on millions of samples, they demonstrate exceptional video

segmentation performance. Recent works improve SAM2 by introducing memory trees [12], motion modeling [27], and memory partitioning [24].

Promptable Concept Detection and Segmentation is formulated as visual grounding, which aims to localize specific targets within an image or video based on user prompts in natural language. Existing approaches mainly fall into two paradigms: the first builds directly upon specialist detection or segmentation architectures, incorporating auxiliary language encoders for precise alignment [17, 18, 23]. The other adopts Large Vision-Language Models (LVLMs) for grounding, either treating grounding as an autoregressive sequence generation task by directly outputting coordinates of bounding boxes [3, 4, 8, 10, 16, 25] or utilizing special tokens to prompt external decoders for fine-grained mask prediction [1, 13, 15, 31]. Despite their success, these approaches struggle to maintain consistent identities in dynamic environments where targets exhibit frequent entries, exits, and reappearances. SAM3 [2] bridges this gap by unifying detection, segmentation, and tracking to ensure robust performance in such dynamic scenarios.

3. Method

3.1. Preliminaries: Memory Selection of SAM3

We focus on simultaneous multi-target video segmentation, in which all objects are inputted to the model at once, rather than inferring for each target at a time. Let $\mathcal{O} = \{o_1, o_2, \dots, o_N\}$ denote targets instantiated in a same initial frame. For each target o_i at time t , SAM3 predicts a mask $M_{i,t}$ and its corresponding segmentation score (or “*query’s score*”) $q_{i,t}$ reflecting the object’s status such as segmentation confidence and individual disappearance status, along with a frame-level presence score p_t demonstrating the overall visibility of all objects in the frame.

Following the design of SAM2 [22], SAM3 does not contain shared object-level contexts and essentially maintains a dedicated memory bank \mathcal{M}_i for each target o_i [2]. Despite this, SAM3 implements a unified memory selection policy across all objects \mathcal{O} , where all N memory banks are either updated together or not at all. Specifically, it calculates an aggregated frame-level confidence score as follows:

$$S_t = \left(\frac{1}{N} \sum_i q_{i,t} \right) \cdot p_t, \quad \forall i \in \{1, \dots, N\}. \quad (1)$$

For target o_i , let $\mathbf{f}_{i,t} = \Phi(I_t, M_{i,t})$ denote the features encoded by the memory encoder Φ and $\mathbf{p}_{i,t}$ be the associated object pointer. The update of its memory bank \mathcal{M}_i at frame t is conditioned on the frame-level average S_t :

$$\mathcal{M}_i \leftarrow \begin{cases} \mathcal{M}_i \cup \{(\mathbf{f}_{i,t}, \mathbf{p}_{i,t})\}, & \text{if } S_t > \tau \\ \mathcal{M}_i, & \text{otherwise} \end{cases} \quad (2)$$

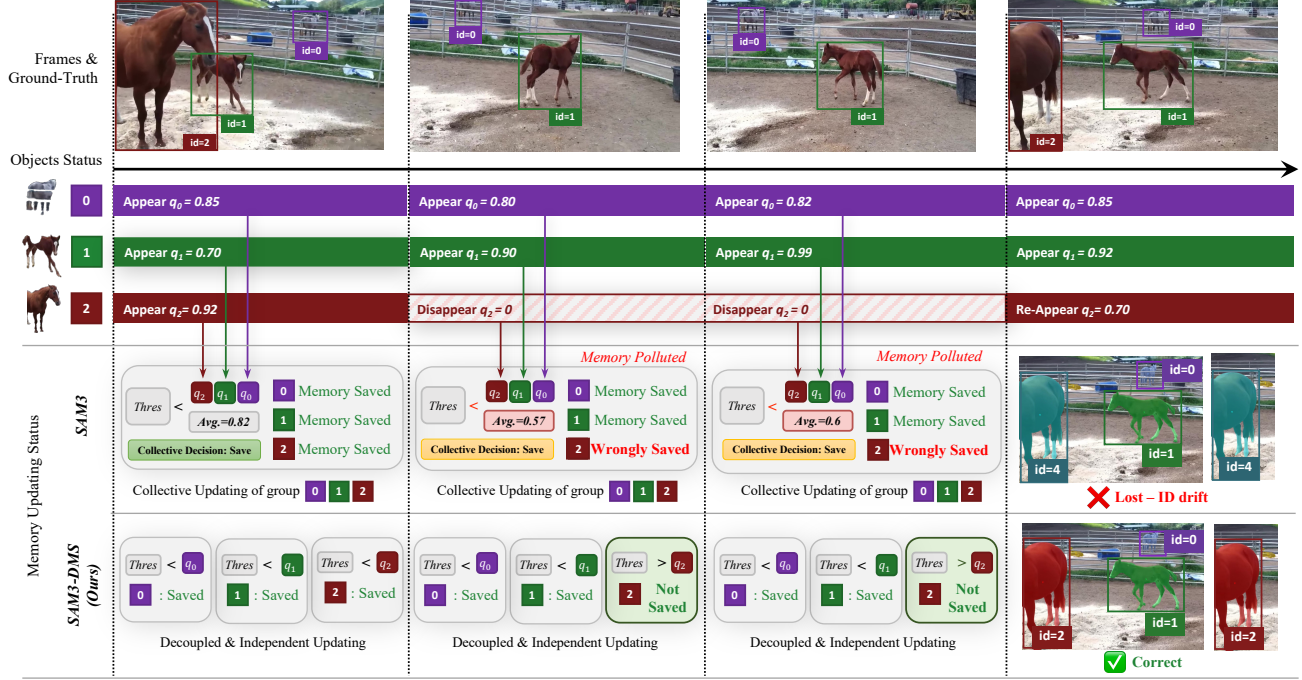


Figure 2. Overview of the Decoupled Memory Selection (DMS) mechanism. In SAM3, the memory status of the group is determined *collectively* by the average score (*Avg*), causing out-of-view objects to be “Wrongly Saved” when other group members remain visible. Our SAM3-DMS evaluates each target’s status *independently*. This decoupling prevents corrupted features from entering the memory bank, leading to the correct identity preservation of Object 2 seen in the final frame.

where τ is a pre-defined confidence threshold. Relying exclusively on this group average, this strategy allows high-performing salient objects to mask the unreliability of uncertain or disappearing objects. As illustrated in Fig. 2, due to the existence of two high-score objects (#0, #1), the overall confidence score exceeds the threshold and causing the model to save all objects’ features, including the disappeared object (#2). This pollutes its memory bank and leads to the following identity drift when the disappeared object re-enters the scene, misidentifying it as a new object (#4).

Notably, for PCS tasks with language prompts where new objects may emerge mid-video, objects that are initialized at the same timestamp are grouped together and all subject to this memory selection policy.

3.2. SAM3-DMS: Decoupled Memory Selection

To mitigate group-level interference and ensure target-level discriminability, we propose the Decoupled Memory Selection (DMS) strategy. Instead of relying on the frame-level S_t , we perform memory selection in an instance-wise manner. This ensures the memory maintenance strictly follows the individual tracking status of its corresponding target. Specifically, the confidence $S_{i,t}$ is computed for each target o_i by combining its own segmentation score $q_{i,t}$ with

the presence score p_t as follows:

$$S_{i,t} = q_{i,t} \cdot p_t, \quad \forall i \in \{1, \dots, N\}, \quad (3)$$

which is then thresholded for the update decision of its corresponding memory bank \mathcal{M}_i :

$$\mathcal{M}_i \leftarrow \begin{cases} \mathcal{M}_i \cup \{(\mathbf{f}_{i,t}, \mathbf{p}_{i,t})\}, & \text{if } S_{i,t} > \tau \\ \mathcal{M}_i, & \text{otherwise} \end{cases} \quad (4)$$

thereby ensuring that memory maintenance is driven exclusively by the quality of each individual target. As shown in Fig. 2, our decoupled approach recognizes the Object #2’s temporary disappearance and excludes blank frames from its memory. This prevents memory pollution, ensuring clean representations for seamless re-identification upon its re-entry in the last frame.

Fundamentally, this decoupled paradigm enables individualized memory update frequencies, allowing each target’s representation to evolve independently according to its specific tracking status. Moreover, since the underlying memory allocation and architecture remain identical to that of SAM3, these improvements are achieved with virtually no additional GPU memory overhead.

Table 1. Quantitative comparisons of PCS on the SA-Co/VEval benchmark. *: Reproduced using official implementation.

Splits	Methods	cgF1	pHOTA	pDetA	pAssA
SA-V val	SAM3	29.3	60.7	44.7	83.2
	SAM3*	29.2	60.68	44.63	83.32
	SAM3-DMS*	29.4 ($\uparrow 0.2$)	60.92 ($\uparrow 0.24$)	44.89 ($\uparrow 0.26$)	83.50 ($\uparrow 0.18$)
SA-V test	SAM3	30.3	58.0	40.9	83.4
	SAM3*	30.1	57.80	40.64	83.30
	SAM3-DMS*	30.3 ($\uparrow 0.2$)	57.97 ($\uparrow 0.17$)	40.83 ($\uparrow 0.19$)	83.41 ($\uparrow 0.11$)
YT-Temporal-1B val	SAM3	50.2	70.5	60.5	82.7
	SAM3*	49.8	70.03	59.64	82.88
	SAM3-DMS*	50.3 ($\uparrow 0.5$)	70.28 ($\uparrow 0.25$)	59.98 ($\uparrow 0.34$)	82.99 ($\uparrow 0.11$)
YT-Temporal-1B test	SAM3	50.8	69.9	60.2	81.7
	SAM3*	49.9	69.17	59.24	81.39
	SAM3-DMS*	51.0 ($\uparrow 1.1$)	69.88 ($\uparrow 0.71$)	59.91 ($\uparrow 0.67$)	82.15 ($\uparrow 0.76$)
SmartGlasses val	SAM3	33.5	60.2	46.2	79.3
	SAM3*	33.2	60.08	46.08	79.02
	SAM3-DMS*	33.6 ($\uparrow 0.4$)	60.29 ($\uparrow 0.21$)	46.17 ($\uparrow 0.09$)	79.41 ($\uparrow 0.39$)
SmartGlasses test	SAM3	36.4	63.6	50.0	81.5
	SAM3*	36.1	63.44	49.90	81.28
	SAM3-DMS*	36.5 ($\uparrow 0.4$)	63.73 ($\uparrow 0.29$)	50.11 ($\uparrow 0.21$)	81.68 ($\uparrow 0.40$)

Table 2. Ablation study on SA-Co/VEval across varying target densities. YT: YT-Temporal-1B, SG: SmartGlasses.

Splits	Methods	≥ 3 targets		≥ 8 targets		≥ 10 targets	
		cgF1	pHOTA	cgF1	pHOTA	cgF1	pHOTA
SA-V test	SAM3	39.78	53.13	41.29	56.30	41.00	53.74
	SAM3-DMS	39.98 ($\uparrow 0.20$)	53.37 ($\uparrow 0.24$)	42.02 ($\uparrow 0.73$)	56.83 ($\uparrow 0.53$)	42.37 ($\uparrow 1.37$)	54.70 ($\uparrow 0.96$)
YT test	SAM3	55.55	67.73	55.82	67.99	57.71	69.64
	SAM3-DMS	56.66 ($\uparrow 1.11$)	68.37 ($\uparrow 0.64$)	57.36 ($\uparrow 1.54$)	68.88 ($\uparrow 0.89$)	59.50 ($\uparrow 1.79$)	70.54 ($\uparrow 0.90$)
SG test	SAM3	48.29	63.48	43.71	61.26	46.29	65.27
	SAM3-DMS	48.80 ($\uparrow 0.51$)	63.75 ($\uparrow 0.27$)	44.46 ($\uparrow 0.75$)	61.84 ($\uparrow 0.58$)	46.81 ($\uparrow 0.52$)	65.65 ($\uparrow 0.38$)

4. Experiments

4.1. Implementation Details

Building upon SAM3, SAM3-DMS is training-free, fully preserving the original model’s PCS and PVS capabilities. All experiments were conducted based on the official implementation of SAM3, on a single NVIDIA RTX A6000 GPU with 48GB of memory.

4.2. Benchmarks and Metrics

PCS. Following [2], we evaluate SAM3-DMS on the SA-Co/VEval benchmark containing 10.3K video-NP pairs, and report the cgF1, pHOTA, pDetA, and pAssA metrics. To further assess the generalizability of the model, we eval-

uate on extra public benchmarks, including YTVIS19 [28], YTVIS21 [29], OVIS [21], and BDD100K [30], following the same evaluation metrics and protocols as [2], where category names are used as conceptual prompts.

PVS. We evaluate SAM3-DMS on SA-V (val/test) [22] and MOSEv2 (val) [9, 11], covering rich multi-target and occlusion scenarios. Following SAM3, we report the $\mathcal{J}\&\mathcal{F}$ metric for SA-V and $\mathcal{J}\&\mathcal{F}$ metric for MOSEv2.

4.3. Quantitative Experiments

PCS. We firstly perform the evaluation on the SA-Co/VEval benchmark. As shown in Table 1, SAM3-DMS consistently outperforms SAM3 on all metrics and datasets. Notably, as demonstrated in Table 2, the performance gap becomes

Table 3. Quantitative comparisons of PCS on other public benchmarks. *: Reproduced using official implementation.

Datasets	Methods	Overall	≥ 3 targets	≥ 5 targets	≥ 7 targets
		mAP	mAP	mAP	mAP
YTVIS19 val	SAM3*	58.38	51.13	48.32	73.35
	SAM3-DMS*	58.41 ($\uparrow 0.03$)	51.22 ($\uparrow 0.09$)	48.86 ($\uparrow 0.54$)	76.46 ($\uparrow 3.11$)
YTVIS21 val	SAM3	57.4	-	-	-
	SAM3*	56.85	53.10	50.94	47.39
	SAM3-DMS*	57.03 ($\uparrow 0.18$)	53.36 ($\uparrow 0.26$)	51.60 ($\uparrow 0.66$)	50.05 ($\uparrow 2.66$)
OVIS val	SAM3	60.5	-	-	-
	SAM3*	61.18	61.07	58.27	54.69
	SAM3-DMS*	62.26 ($\uparrow 1.08$)	62.13 ($\uparrow 1.06$)	59.94 ($\uparrow 1.67$)	56.19 ($\uparrow 1.50$)
		TETA	HOTA	DetA	AssA
BDD100K val	SAM3	47.2	-	-	-
	SAM3*	49.49	40.45	31.26	55.13
	SAM3-DMS*	49.71 ($\uparrow 0.22$)	40.75 ($\uparrow 0.30$)	31.50 ($\uparrow 0.24$)	55.52 ($\uparrow 0.39$)

Table 4. Quantitative comparisons of PVS on different benchmarks. *: Reproduced using official implementation.

Methods	Mode	SA-V val			SA-V test			MOSEv2 val		
		$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}
SAM3	One-by-one	83.5	-	-	84.4	-	-	60.3	-	-
SAM3*	One-by-one	83.3	79.4	87.2	84.3	80.4	88.3	60.3	57.9	62.7
SAM3	Simultaneous	81.2	77.4	84.9	81.3	77.4	85.2	60.0	57.6	62.4
SAM3-DMS	Simultaneous	83.3 ($\uparrow 2.1$)	79.4	87.2	84.3 ($\uparrow 3.0$)	80.4	88.3	60.3 ($\uparrow 0.3$)	57.9	62.7

more significant with increasing target density. Specifically, as the scenes become extremely complex and contain more than 10 targets, our performance surpasses the original SAM3 by 1.79 cgF1 on the YT-Temporal-1B test set and 1.37 on the SA-V test set. Moreover, a similar trend is observed on other public benchmarks, as in Table 3. For example, on the YTVIS21 validation set, the improvement of mAP scales from 0.18 (overall) to 0.66 (≥ 5 targets), ultimately reaching 2.66 (≥ 7 targets). These results demonstrate the effectiveness of our proposed DMS strategy in preserving identity for multi-target segmentation.

PVS. We additionally benchmark SAM3-DMS for simultaneous multi-target PVS. As shown in Table 4, compared with the upper bound where each target is independently inputted and inferred by the network one-by-one, the shared memory selection used in the original SAM3 suffers from severe performance degradation. Our approach effectively bridges this gap, achieving gains of 2.1 and 3.0 $\mathcal{J}\&\mathcal{F}$ on SA-V val and test splits, and achieves almost no performance loss on both SA-V and MOSEv2 compared with the upper bound. Note in PCS, one-by-one inference is not applicable due to the dynamic emergence of new objects.

4.4. Qualitative Experiments

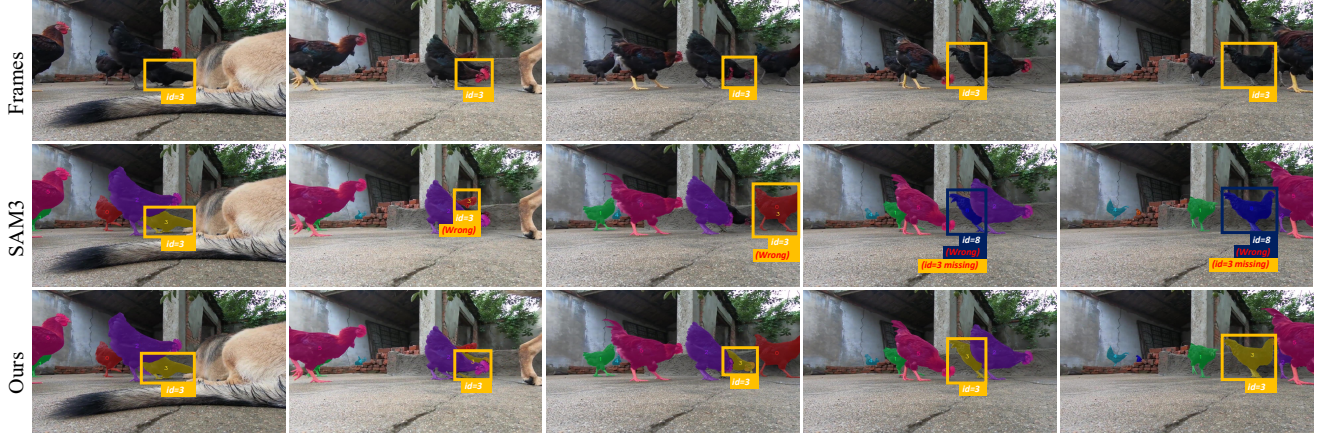
We present qualitative comparisons on PCS against the SAM3 baseline on challenging scenarios in Figs. 3 to 6.

Case 1: ID switch on heavy occlusion.

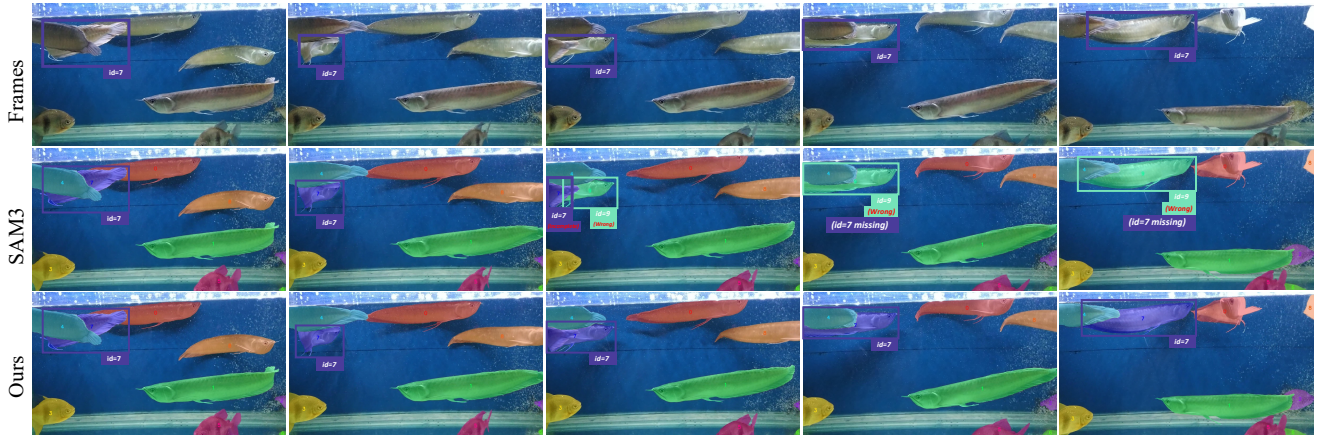
As shown in Fig. 3, when a target is occluded, its visual cues become incomplete, such as chicken neck (top) or fish tail (bottom). These incomplete cues fail to accurately represent the target’s biological features. SAM3 encodes this information into memory based on the average group-level score, leading to subsequent identity drifts. In contrast, SAM3-DMS discards unclear frames, ensuring only clear representations are stored, preventing ID switches.

Case 2: ID switch on disappearance and re-appearance.

As shown in Fig. 4, when a target exits the view, high confidence from other visible objects might cause SAM3 to sustain a high group average, polluting the target’s memory with noisy features, resulting in unrecoverable tracking loss. SAM3-DMS excludes out-of-view frames, ensuring valid and effective object memory for re-identification. We provide additional visualizations in Fig. 4b and 4c, featuring scenarios of multiple targets exiting and re-appearance, highlighting the robustness of our approach.



(a) “The chickens”



(b) “The fishes”

Figure 3. Example results of heavy occlusions of targets. (Case 1, zoom in for better view.)

Case 3: Interference from distractors.

As illustrated in Fig. 5, tracking errors frequently occur when distractors exhibit parallel motion as in the crowded fish groups of Fig. 5a, or obstruct the target’s path as in the pedestrians of Fig. 5b. In these scenarios, distractors tend to overshadow the target and significantly suppress its confidence score, thereby preventing memory updates, while SAM3-DMS maintains independent updates for each target, effectively circumventing distractor interference even in very dense and crowded scene layouts.

Case 4: Random ID drifts.

As shown in Fig. 6, SAM3 suffers from random and unexpected ID drifts, acting as “hallucinations” where false positives are generated on totally unrelated objects. This typically occurs under rapid motion, such as the fast-moving plane in Fig. 6a or the rapidly walking pedestrian in Fig. 6b. Our decoupled approach effectively captures the appearance and motion of each individual target, knowing when a target exits the scene, thereby preventing such drifts.

5. Discussion and Conclusion

Limitations. We aim to improve the multi-target video segmentation robustness by improving the memory updating strategy. We enhance the performance for PCS, however, for PVS, we mitigate the performance gap between the separate one-by-one inference of each target and simultaneous multi-target inference, but without improving its upper bound. Future research could be conducted to elevate the upper bound of SAM3 for PVS.

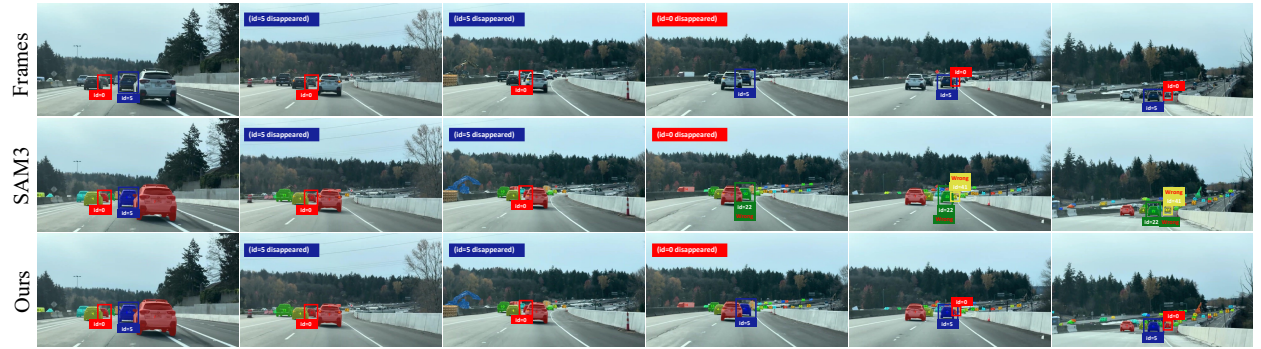
Conclusion. In this work, we propose SAM3-DMS, addressing the limitations of SAM3’s synchronized memory selection in simultaneous multi-target video segmentation. By introducing a training-free decoupled memory selection strategy, we address the memory pollution issues. Experiments demonstrate that our approach enhances multi-target robustness in various challenging benchmarks, with performance gains scaling positively with target density, establishing a robust foundation for video grounding in the wild.



(a) "The ducks"

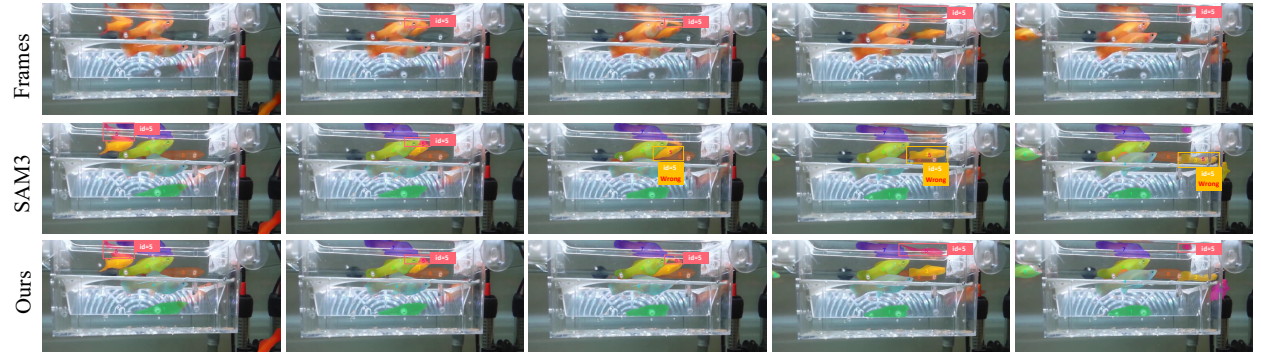


(b) "Men with jeans"



(c) "Vehicles"

Figure 4. Example results of disappearance and re-appearance of targets. (Case 2, zoom in for better view.)



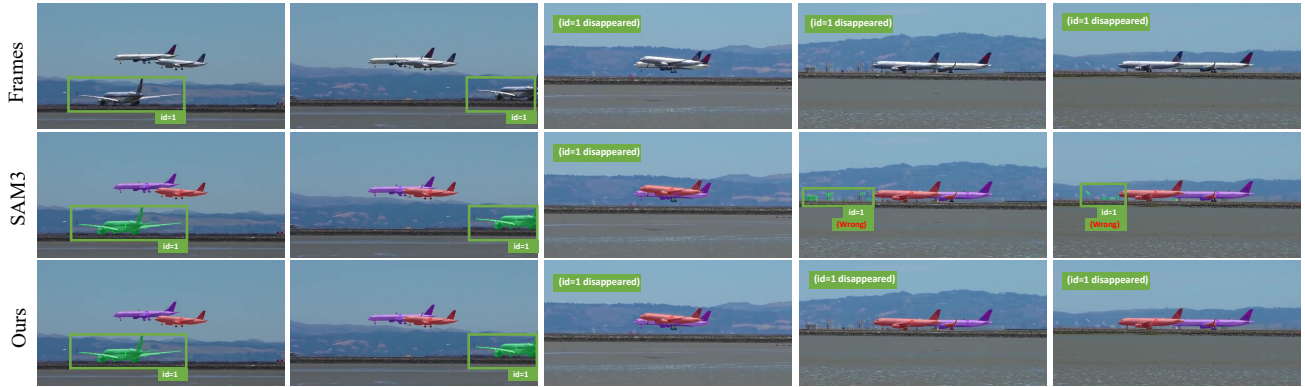
(a) "The fishes"

Figure 5. Example results of object interferences, where distractors exhibit parallel motion. (Case 3, zoom in for better view.)



(b) "Pedestrians"

Figure 5. (cont.) Example results of object interferences, where distractors obstruct the target's path. (Case 3, zoom in for better view.)



(a) "The planes"



(b) "Person with handbag"

Figure 6. Example results of rapid motion, which cause SAM3 to exhibit random ID drifts on unexpected objects. (Case 4, zoom in for better view.)

References

- [1] Zechen Bai, Tong He, Haiyang Mei, Pichao Wang, Ziteng Gao, Joya Chen, Zheng Zhang, and Mike Zheng Shou. One token to seg them all: Language instructed reasoning segmentation in videos. *Advances in Neural Information Processing Systems*, 37:6833–6859, 2024. 2
- [2] Nicolas Carion, Laura Gustafson, Yuan-Ting Hu, Shoubhik Debnath, Ronghang Hu, Didac Suris, Chaitanya Ryali, Kalyan Vasudev Alwala, Haitham Khedr, Andrew Huang, et al. Sam 3: Segment anything with concepts. *arXiv preprint arXiv:2511.16719*, 2025. 1, 2, 4
- [3] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multi-modal llm’s referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023. 2
- [4] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024. 2
- [5] Ho Kei Cheng and Alexander G Schwing. Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model. In *European conference on computer vision*, pages 640–658. Springer, 2022. 2
- [6] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Rethinking space-time networks with improved memory coverage for efficient video object segmentation. *Advances in neural information processing systems*, 34:11781–11794, 2021. 2
- [7] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Joon-Young Lee, and Alexander Schwing. Putting the object back into video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3151–3161, 2024. 2
- [8] Henghui Ding, Chang Liu, Shuting He, Xudong Jiang, and Chen Change Loy. MeViS: A large-scale benchmark for video segmentation with motion expressions. In *ICCV*, 2023. 2
- [9] Henghui Ding, Chang Liu, Shuting He, Xudong Jiang, Philip HS Torr, and Song Bai. MOSE: A new dataset for video object segmentation in complex scenes. In *ICCV*, 2023. 4
- [10] Henghui Ding, Chang Liu, Shuting He, Kaining Ying, Xudong Jiang, Chen Change Loy, and Yu-Gang Jiang. MeViS: A multi-modal dataset for referring motion expression video segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025. 2
- [11] Henghui Ding, Kaining Ying, Chang Liu, Shuting He, Xudong Jiang, Yu-Gang Jiang, Philip HS Torr, and Song Bai. MOSEv2: A more challenging dataset for video object segmentation in complex scenes. *arXiv preprint arXiv:2508.05630*, 2025. 4
- [12] Shuangrui Ding, Rui Qian, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Yuwei Guo, Dahua Lin, and Jiaqi Wang. Sam2long: Enhancing sam 2 for long video segmentation with a training-free memory tree. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13614–13624, 2025. 2
- [13] Sitong Gong, Yunzhi Zhuge, Lu Zhang, Zongxin Yang, Pingping Zhang, and Huchuan Lu. The devil is in temporal token: High quality video reasoning segmentation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 29183–29192, 2025. 2
- [14] 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 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 1
- [15] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589, 2024. 2
- [16] Hongyu Li, Jinyu Chen, Ziyu Wei, Shaofei Huang, Tianrui Hui, Jialin Gao, Xiaoming Wei, and Si Liu. Llava-st: A multimodal large language model for fine-grained spatial-temporal understanding. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 8592–8603, 2025. 2
- [17] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10965–10975, 2022. 2
- [18] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European conference on computer vision*, pages 38–55. Springer, 2024. 2
- [19] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9226–9235, 2019. 2
- [20] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alex Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv preprint arXiv:1704.00675*, 2017. 2
- [21] Jiyang Qi, Yan Gao, Yao Hu, Xinggang Wang, Xiaoyu Liu, Xiang Bai, Serge Belongie, Alan Yuille, Philip HS Torr, and Song Bai. Occluded video instance segmentation: A benchmark. *International Journal of Computer Vision*, 130(8): 2022–2039, 2022. 4
- [22] 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. 1, 2, 4
- [23] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, et al. Grounded sam: Assembling open-world models for diverse visual tasks. *arXiv preprint arXiv:2401.14159*, 2024. 2
- [24] Jovana Videnovic, Alan Lukezic, and Matej Kristan. A distractor-aware memory for visual object tracking with

- sam2. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24255–24264, 2025. [2](#)
- [25] Jiankang Wang, Zhihan Zhang, Zhihang Liu, Yang Li, Jian-nan Ge, Hongtao Xie, and Yongdong Zhang. Spacevlm: Endowing multimodal large language model with spatio-temporal video grounding capability. *arXiv preprint arXiv:2503.13983*, 2025. [2](#)
- [26] Haozhe Xie, Hongxun Yao, Shangchen Zhou, Shengping Zhang, and Wenxiu Sun. Efficient regional memory network for video object segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1286–1295, 2021. [2](#)
- [27] Cheng-Yen Yang, Hsiang-Wei Huang, Wenhao Chai, Zhongyu Jiang, and Jenq-Neng Hwang. Samurai: Adapting segment anything model for zero-shot visual tracking with motion-aware memory. *arXiv preprint arXiv:2411.11922*, 2024. [2](#)
- [28] Linjie Yang, Yuchen Fan, and Ning Xu. Video instance segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5188–5197, 2019. [4](#)
- [29] Linjie Yang, Yuchen Fan, Yang Fu, and Ning Xu. The 3rd large-scale video object segmentation challenge-video instance segmentation track. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshop*, 2021. [4](#)
- [30] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madhavan, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2636–2645, 2020. [4](#)
- [31] Haobo Yuan, Xiangtai Li, Tao Zhang, Zilong Huang, Shilin Xu, Shunping Ji, Yunhai Tong, Lu Qi, Jiashi Feng, and Ming-Hsuan Yang. Sa2va: Marrying sam2 with llava for dense grounded understanding of images and videos. *arXiv preprint arXiv:2501.04001*, 2025. [2](#)