# MASSeg : 2nd Technical Report for 4th PVUW MOSE Track

Xuqiang Cao[1]    Linnan Zhao[1]    Jiaxuan Zhao[1]    Fang Liu[1]
Puhua Chen[1*]    Wenping Ma[1]

[1]Key Laboratory of Intelligent Perception and Image Understanding
Xi'an, China

## Abstract

*Complex video object segmentation continues to face significant challenges in small object recognition, occlusion handling, and dynamic scene modeling. This report presents our solution, which ranked **second** in the MOSE track of CVPR 2025 PVUW Challenge. Based on an existing segmentation framework, we propose an improved model named **MASSeg** for complex video object segmentation, and construct an enhanced dataset, **MOSE+**, which includes typical scenarios with occlusions, cluttered backgrounds, and small target instances. During training, we incorporate a combination of inter-frame consistent and inconsistent data augmentation strategies to improve robustness and generalization. During inference, we design a mask output scaling strategy to better adapt to varying object sizes and occlusion levels. As a result, MASSeg achieves a **J score of 0.8250**, **F score of 0.9007**, and a **J&F score of 0.8628** on the MOSE test set. The code is available at https://github.com/cxqNet/MASSeg.*

## 1. Introduction

Pixel-level scene understanding is a fundamental task in computer vision, aiming to recognize the category, semantics, and instance of each pixel. With the rapid growth of video content, this task has extended from static images to dynamic videos, drawing increasing attention to video object segmentation (VOS). VOS requires segmenting and tracking target objects across frames with only the first-frame mask provided. It has been widely applied in autonomous driving, augmented reality, video editing, and data annotation.

As task complexity increases, memory-based VOS methods have become the mainstream [3]. These models maintain a memory bank to store visual features and adopt attention or feature matching to segment targets in future frames. Representative works include STM [7], which introduces explicit memory for pixel-level matching; STCN [12],



Figure 1. Representative challenges in complex video object segmentation. The examples showcase small object motion (left), appearance confusion with occlusion (middle), and densely cluttered scenes (right), which reflect typical situations encountered in the MOSE dataset.

which improves efficiency with unmasked frames and L2 affinity; and XMem [2], which mimics human memory with sensory, working, and long-term stages. Cutie [11] enhances robustness in crowded scenes via object-level memory and object transformers, achieving state-of-the-art results on YouTube-VOS and DAVIS.

To better evaluate model robustness in real-world scenarios, Henghui Ding et al. proposed and organized the MOSE [6] and MeViS [4] tracks. These tracks target two distinct but complementary aspects of video understanding: MOSE focuses on complex video object segmentation under challenging visual conditions, while MeViS emphasizes motion expression-guided segmentation based on natural language descriptions.

To better evaluate model robustness in real-world scenarios, Henghui Ding et al., in collaboration with CVPR, introduced the Pixel-level Video Understanding in the Wild (PVUW) Challenge. The challenge focuses on pixel-level scene understanding and features two complementary tracks: MOSE [5], which targets complex multi-object segmentation under challenging visual conditions, and MeViS [4], which emphasizes motion expression-guided segmentation based on natural language descriptions. The MOSE track achieved significant progress in the PVUW 2024 Challenge [6].

The MOSE dataset includes 2,149 videos and 5,200 annotated instances over 36 categories, with over 430K pixel-level masks. MOSE emphasizes real-world challenges such
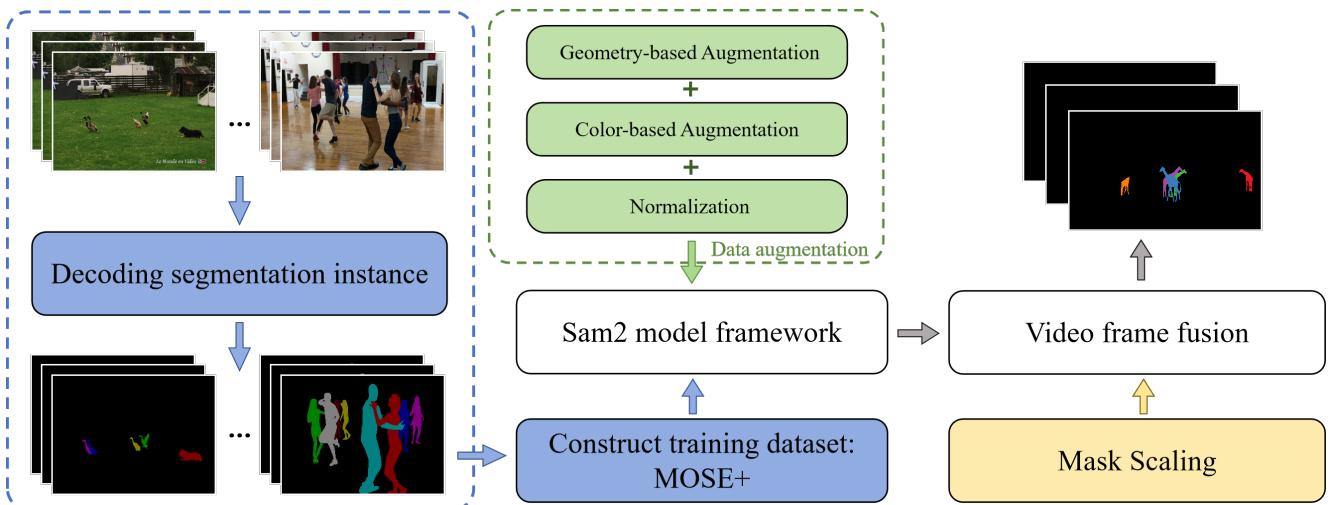
Figure 2. Overview of our method.

as occlusions, fast motion, dense small objects, similar appearances, and frequent reappearance, providing a rigorous benchmark beyond prior datasets like DAVIS and YouTube-VOS. As illustrated in Figure 1, the MOSE dataset contains challenging scenarios such as small objects, heavy occlusions, and similar instances, which significantly impact segmentation robustness.

Mainstream methods, e.g., Cutie, achieve over 84% J&F on YouTube-VOS but drop significantly on MOSE, revealing generalization limitations. MOSE has thus emerged as a valuable platform for driving progress in robust, real-world VOS.

To address MOSE-specific challenges, we propose an improved method based on the SAM2 framework [10], incorporating a mask output control strategy and retraining on an enhanced dataset **MOSE+**, which includes diverse small object instances and complex scenes. We further adopt inter-frame *consistent* and *inconsistent* augmentations and tuned inference strategies to improve robustness across scales and occlusions. Our method achieves a **J score of 0.8250**, **F score of 0.9007**, and a **J&F of 0.8628** on the MOSE test set, ranking 2nd overall.

## 2. Method

Our overall solution is illustrated in Figure 2. Based on the data characteristics of the MOSE dataset, we construct a customized dataset named **MOSE+** and design a series of targeted data augmentation strategies, aiming to simulate appearance variations, pose perturbations, illumination changes, and structural blur commonly found in real-world video scenarios. During inference, a *mask confidence control strategy* is adopted, followed by a fusion process over video frames to obtain the final segmentation results. The detailed components are described below.

### 2.1. Baseline Model

We adopt a transformer-based segmentation framework featuring object-guided attention, mask-aware memory, and spatiotemporal reasoning. The model effectively captures temporal cues and spatial details through dual memory modules and multi-scale decoding, enabling robust performance under challenging scenarios such as occlusion, motion blur, and small-object clutter. This strong baseline lays a solid foundation for our enhancement strategies.

### 2.2. Loss Function

To achieve high-precision segmentation and temporal consistency, we design a multi-task loss that combines pixel-wise accuracy, region-level overlap, classification discriminability, and robustness to occlusion. The total loss is defined as:

$$\mathcal{L}_{total} = \lambda_1 \mathcal{L}_{CE} + \lambda_2 \mathcal{L}_{Dice} + \lambda_3 \mathcal{L}_{Sim} + \lambda_4 \mathcal{L}_{MaskIoU} \quad (1)$$

where $\mathcal{L}_{CE}$ denotes cross-entropy loss for foreground-background classification, $\mathcal{L}_{Dice}$ enhances region consistency, $\mathcal{L}_{Sim}$ enforces similarity between memory and query features, and $\mathcal{L}_{MaskIoU}$ constrains predicted mask quality. These losses are computed across multiple frames and candidate masks to jointly supervise spatiotemporal modeling.

In our implementation, the weights $\lambda_i$ are empirically set for balanced optimization. This formulation improves segmentation performance on small or occluded targets and supports temporal consistency in long sequences.

### 2.3. Data Augmentation

To improve generalization and robustness, we introduce a set of targeted augmentation strategies during training.
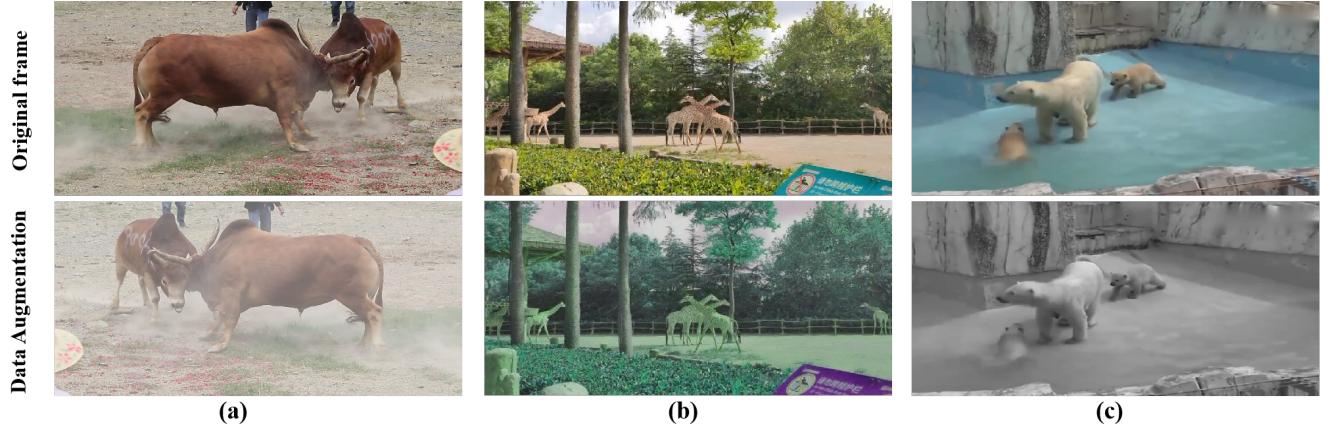
Figure 3. Visualization of our data augmentation strategies. Each column shows an example before (top) and after (bottom) applying augmentation. From left to right: (a) geometric transformation (e.g., affine distortion), (b) color jittering with inconsistent color shift, and (c) grayscale conversion. These augmentations simulate realistic variations in pose, illumination, and appearance, improving robustness and generalization of the model in complex scenarios.

Unlike static image tasks, video segmentation demands consistency across frames while simulating realistic variations. Our approach integrates both frame-consistent and frame-inconsistent perturbations:

- **Consistent geometric transformations**: Random horizontal flipping, affine transformations (rotation, shear), and multi-scale resizing are applied across all frames in a clip to simulate viewpoint and object deformation.
- **Mixed color perturbations**: Brightness, contrast, and saturation changes are applied globally, while grayscale conversion and inconsistent color jittering are selectively applied to individual frames, enhancing robustness to lighting changes and visual ambiguity.
- **Normalization**: Images are transformed into tensors and normalized using ImageNet mean and standard deviation for stable convergence and pretrained compatibility.

These augmentations significantly improve the model's ability to handle structure variation, appearance change, and dynamic scenes in MOSE-like scenarios.

### 2.4. Inference Strategy

To enhance the robustness and adaptability of the model in complex video object segmentation, we design and apply a series of strategies during inference.

**Mask Confidence Control Strategy.** We observe that the quality of predicted masks can be significantly affected by post-processing in different scenarios, such as small objects, heavy occlusions, and target overlaps. To address this, we adopt a control strategy based on dynamic adjustment of the mask output distribution, using two key parameters: *sigmoid scale* and *sigmoid bias*. The sigmoid scale controls the sharpness of the output boundaries, while the sigmoid bias adjusts the overall activation level, thereby influencing the

target coverage and boundary quality.

Extensive experiments on the validation set show that setting the sigmoid scale to 7.5 and the sigmoid bias to -4.0 yields the best performance on the MOSE dataset in terms of J&F metrics. In addition, other parameter combinations are explored to improve local segmentation performance under specific conditions such as dense small objects or complex backgrounds. Final predictions are obtained by merging results from different local configurations to achieve globally optimal segmentation.

### 3. Implementation Details

**Data.** To improve generalization and target modeling in complex scenarios, we construct an enhanced training set named **MOSE+**, based on the original MOSE dataset. This augmented set is composed of video segments from multiple public VOS datasets, selected to match the characteristics of MOSE, including frequent occlusions, dense small objects, object reappearance, and high similarity among targets. Specifically, we integrate carefully chosen sequences from datasets such as BURST [1], DAVIS [8], OVIS [9], and YouTubeVIS [13], unify their annotations and resolution formats, and seamlessly merge them with MOSE to form a consistent training set, enhancing semantic understanding and robustness.

**Training.** We train our model end-to-end on the MOSE+ dataset. Each batch contains 2 samples, and we adopt the AdamW optimizer with an initial learning rate of 1e-5. The input image resolution is $1024 \times 1024$, with each clip containing 6 frames and supporting up to 10 object instances. We employ automatic mixed precision (AMP) and apply gradient clipping (max norm = 0.1) to stabilize convergence. The training is conducted on three NVIDIA H100
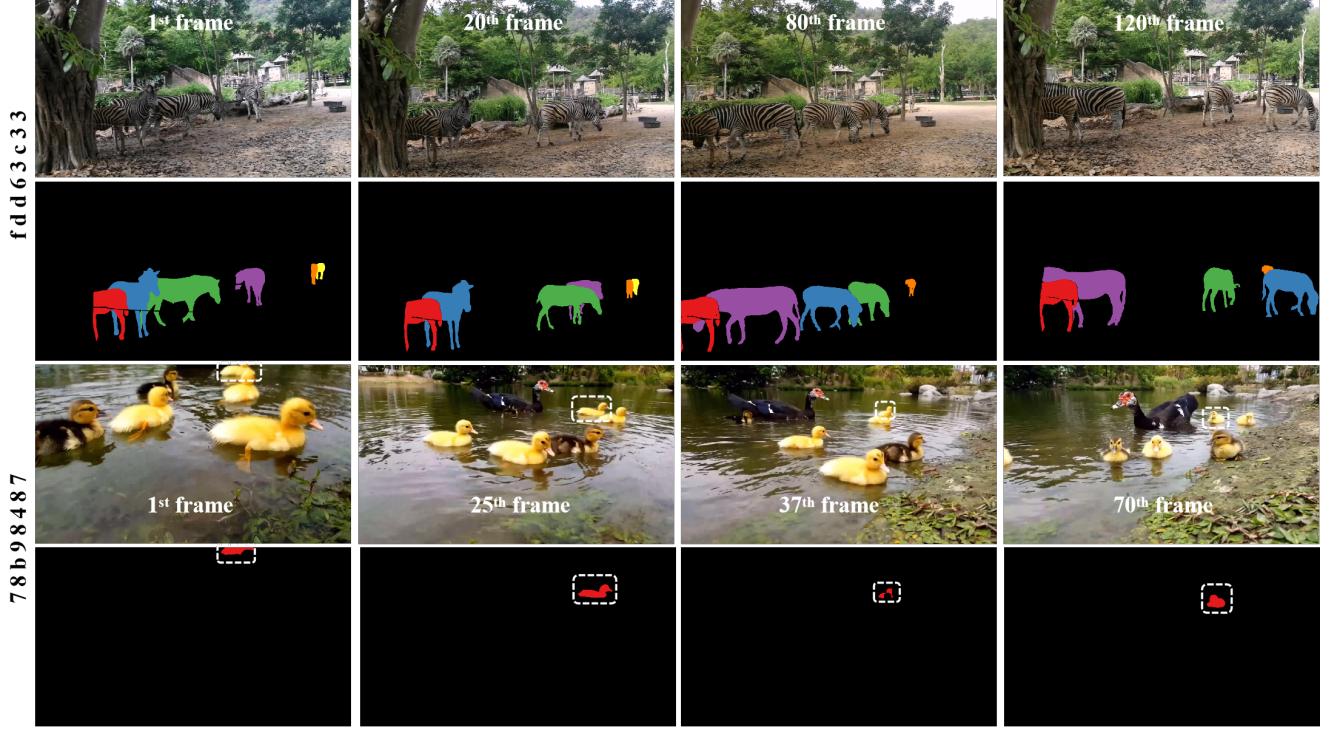
Figure 4. Qualitative results of our method on challenging MOSE test sequences. Our model accurately segments small objects, handles severe occlusions, and maintains temporal consistency across fast-moving and cluttered scenes.

NVLink GPUs (each with 96GB memory) for 50 epochs, taking approximately 40 hours in total.

### 3.1. Main Results

| Method | J&F | J | F |
| --- | --- | --- | --- |
| imaplus (1st) | 0.8726 | 0.8359 | 0.9092 |
| KirinCZW (ours) | **0.8628** | **0.8250** | **0.9007** |
| dumplings (3rd) | 0.8392 | 0.8028 | 0.8757 |

Table 1. Leaderboard of the MOSE Track in the 4th PVUW Challenge 2025.

Our method achieved second place in the MOSE Track of the CVPR 2025 PVUW Challenge, with a J&F score of **0.8628**, consisting of a region similarity (J) of **0.8250** and a contour accuracy (F) of **0.9007**. The leaderboard is summarized in Tab. 1, demonstrating the robustness and effectiveness of our method in complex video object segmentation scenarios.

### 3.2. Ablation Study

To evaluate the contribution of each component, we conduct ablation studies on the MOSE validation set. As shown in Tab. 2, we start from the public baseline Cutie, which achieved a J&F of 0.7065 on MOSE—significantly lower than its performance on traditional datasets like YouTube-VOS, indicating the challenges posed by complex scenes.

| Method | J | F | J&F |
| --- | --- | --- | --- |
| Cutie (val) | 0.6511 | 0.7619 | 0.7065 |
| Baseline (val) | 0.6953 | 0.7761 | 0.7357 |
| Baseline + DA (val) | 0.7181 | 0.7947 | 0.7564 |
| Baseline + DA + MSS (val) | 0.7339 | 0.8191 | 0.7765 |
| Baseline + DA + MSS (test) | **0.8250** | **0.9007** | **0.8628** |

Table 2. Ablation study results on the MOSE validation and test sets.

We then evaluate our baseline framework with multi-scale encoders, enhanced memory, and mask mechanisms, which improves J&F to 0.7357. Adding MOSE+ and our proposed data augmentation (DA) strategy further raises performance to 0.7564, verifying the effectiveness of enhanced training data. Finally, we introduce our Mask Scaling Strategy (MSS) to dynamically adjust the output mask distribution using *sigmoid scale* and *bias*, which boosts the performance to 0.7765 on the validation set and **0.8628** on the test set. As shown by the qualitative results in Fig. 4, our optimized method demonstrates superior performance on the MOSE test set, particularly in segmenting small objects, occluded targets, and objects within complex and clut-

tered scenes.

## 4. Conclusion

In this work, we proposed a robust solution for complex video object segmentation by integrating enhanced training strategies and a mask confidence control mechanism. Based on the MOSE+ dataset, our approach incorporates targeted data augmentation and adaptive mask output calibration to improve segmentation performance under challenging scenarios. Our method achieved a J&F score of 0.8628 and ranked second in the MOSE track of the CVPR 2025 PVUW Challenge.

## References

[1] Ali Athar, Jonathon Luiten, Paul Voigtlaender, Tarasha Khurana, Achal Dave, Bastian Leibe, and Deva Ramanan. Burst: A benchmark for unifying object recognition, segmentation and tracking in video. In *WACV*, 2023. 3

[2] Yuxin Cheng, Yihao Liu, Zeyu Wang, Ke Li, Yu-Wing Tai, and Chi-Keung Tang. Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model. In *ECCV*, 2022. 1

[3] Henghui Ding, Xudong Jiang, Bing Shuai, Ai Qun Liu, and Gang Wang. Semantic segmentation with context encoding and multi-path decoding. *IEEE Transactions on Image Processing*, 2020. 1

[4] Henghui Ding, Chang Liu, Shuting He, Xudong Jiang, and Chen Change Loy. Mevis: A large-scale benchmark for video segmentation with motion expressions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2694–2703, October 2023. 1

[5] 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. 1

[6] Henghui Ding, Chang Liu, Yunchao Wei, Nikhila Ravi, Shuting He, Song Bai, Philip Torr, Deshui Miao, Xin Li, Zhenyu He, Yaowei Wang, Ming-Hsuan Yang, Zhensong Xu, Jiangtao Yao, Chengjing Wu, Ting Liu, Luoqi Liu, Xinyu Liu, Jing Zhang, Kexin Zhang, Yuting Yang, Licheng Jiao, Shuyuan Yang, Mingqi Gao, Jingnan Luo, Jinyu Yang, Jungong Han, Feng Zheng, Bin Cao, Yisi Zhang, Xuanxu Lin, Xingjian He, Bo Zhao, Jing Liu, Feiyu Pan, Hao Fang, and Xiankai Lu. Pvuw 2024 challenge on complex video understanding: Methods and results, 2024. 1

[7] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *ICCV*, 2019. 1

[8] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alexander Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv:1704.00675*, 2017. 3

[9] Jiyang Qi, Yan Gao, Yao Hu, Xinggang Wang, Xiaoyu Liu, Xiang Bai, Serge Belongie, Alan Yuille, Philip Torr, and Song Bai. Occluded video instance segmentation: A benchmark. *International Journal of Computer Vision*, 2022. 3

[10] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 2

[11] Tao Tang, Xiaoyang Wu, Zhihao Chen, Xinggang Wang, Wenyu Liu, and Xiang Bai. Associating objects with transformers for video object segmentation. In *CVPR*, 2022. 1

[12] Angtian Wang, Linjie Yang, Zhe Lin, Kevin Barnes, and Humphrey Shi. Rethinking space-time networks with improved memory coverage for efficient video object segmentation. In *NeurIPS*, 2021. 1

[13] Linjie Yang, Yuchen Fan, and Ning Xu. The 2nd large-scale video object segmentation challenge - video object segmentation track, Oct. 2019. 3