# Benchmarking SAM2-based Trackers on FMOX

Senem Aktas[0000−0002−3996−2771], Charles Markham[0000−0003−2447−2611], John
McDonald[0000−0001−9225−673X], and Rozenn Dahyot[0000−0003−0983−3052]

Department of Computer Science, Maynooth University, Ireland

**Abstract.** Several object tracking pipelines extending Segment Any-
thing Model 2 (SAM2) have been proposed in the past year, where
the approach is to follow and segment the object from a single exem-
plar template provided by the user on a initialization frame. We pro-
pose to benchmark these high performing trackers (SAM2, EfficientTAM,
DAM4SAM and SAMURAI) on datasets containing fast moving objects
(FMO) specifically designed to be challenging for tracking approaches.
The goal is to understand better current limitations in state-of-the-art
trackers by providing more detailed insights on the behavior of these
trackers. We show that overall the trackers DAM4SAM and SAMURAI
perform well on more challenging sequences.

**Keywords:** Segment anything model · Fast moving object · Tracking.

## 1 Introduction

Video Object Tracking (VOT) and Video Object Segmentation (VOS) both aim
to follow the object throughout a video or image sequence. VOT, focuses on
locating and following the position of a target object, typically outputting a
bounding box without providing detailed shape or pixel-level information. While
VOS, seeks to identify and segment the object at the pixel level, producing a
mask in each frame, thereby capturing its shape and boundaries [3].

SAM2 is one of the state-of-the-art (SOTA) VOS/VOT methods that uses
object positions given as points, bounding boxes, or masks in any frame to
initialize tracking [3]. Various extentions of SAM2 has been proposed for specific
purposes; for instance DAM4SAM [11] for handling distractors, SAMURAI [13]
managing fast motions, and EfficientTAM [12] for improving efficiency across
various platforms (cf. Section 2.2).

Recent SAM2-based works proposed the DiDi dataset [11] to address distrac-
tors, and the Mosev2 dataset [3] to handle various challenging cases. The study
in [3] utilized SAM2 and its variants, including DAM4SAM [11] and SAMURAI
[13]. Similarly, Aktas et al [1] introduced FMOX, a JSON format designed for
challenging Fast Moving Object (FMO) datasets, and extended the ground truth
annotations to include object size categorization. FMOX [1] has been used to
evaluate the SAM-based tracker EfficientTAM, demonstrating its performance
compared to FMO-specific pipelines [10,5,9,6,8] using the Trajectory Intersection
over Union (TIoU) metric.

In this paper, we extend the benchmarking of SAM2 and its variants DAM4SAM, and SAMURAI, alongside EfficientTAM on FMOX dataset through the use of more standard performance metrics: the Mean Intersection over Union (mIoU) and the Dice Score metrics. These generalized metrics were prioritized to facilitate a broader comparison against the wider state-of-the-art literature. The results indicate that DAM4SAM consistently outperforms the other trackers on FMO datasets, aligning with similar observations from Mosev2 [3]. On the other hand, EfficientTAM shows comparatively lower performance among the SAM2-based trackers examined. In Section 2, we provide a brief background on the datasets and SAM2-based trackers. Section 3 presents our methodology for benchmarking along with the tracker initialisation process and performance analysis. In Section 4, we present and discuss our findings, and finally, Section 5 summarises the conclusions drawn from this study.

## 2   Background

### 2.1   Benchmark datasets for object tracking

Despite the abundance of video and image benchmarks, many challenges remain unaddressed in object tracking which limit their ability to generalize to complex real-world scenarios [3]. The recently proposed MOSEv2 dataset (coMplex video Object SEgmentation) [3] addresses several of these difficulties by including videos with complex scenes featuring object disappearance and reappearance, heavy occlusions, crowded areas, small objects, poor lighting, and camouflage. Ding et al [3] use MOSEv2 dataset to evaluate trackers such as SAM2 [7], SAMURAI [13], and DAM4SAM [11].

Even though some challenges remain in standard benchmarks, the overall results reported often fail to reflect these difficulties, as many trackers do not effectively capture or address them. This leads to inflated performance scores that mask the true complexity of real-world tracking scenarios [11]. Addressing this gap, DAM4SAM [11] focused on distractors and occlusions by carefully selecting validation and test sequences from major benchmarks including LaSOT [4] and GOT-10k [2], forming the DiDi dataset to enable more rigorous evaluation.

A similar argument can be made for Fast Moving Objects (FMOs), which represent a significant yet often overlooked challenge in tracking due to their high speed and motion-induced blur. These characteristics limit the effectiveness of many current trackers and are not adequately captured by standard benchmarks. This is particularly important given that the FMO problem is essential for advancing tracking performance in practical applications, such as sports analysis, autonomous driving, and robotics.

We focus here on datasets specifically designed for evaluating FMOs, including Falling Object (6 sequences, [5]), TbD (12 sequences, [6]), TbD-3D (10 sequences, [8]) and FMOv2 for which 18 sequences are used here from the 19 available ([9], the sequence `more_balls` with multiple objects has been excluded from our analysis due to its multi-instance objects lacking unique IDs [1]). Each

dataset offers unique features: for instance, TbD-3D extends the challenge by incorporating 3D object motion and appearance changes. While FMOv2 is designed to be more challenging, with high object displacements and almost no bounding box overlap (IoU) between consecutive frames. These datasets (collectively named the FMOX dataset) have been recently augmented with ground truth JSON files to enable straightforward and easy-to-use benchmarking of trackers [1]. FMO datasets not only include fast-moving objects but also small objects where the FMOX description can be used to focus on specific object size categories [1] such as small objects which are challenging for tracking (c.f. Fig. 1). We provide benchmark results on the FMOX dataset in Section 4 for showcasing the capabilities of these various trackers.
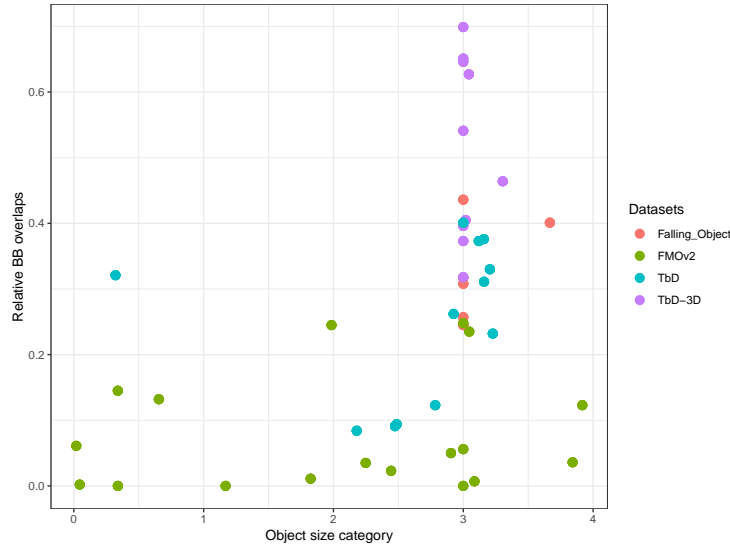


**Fig. 1.** FMOX regroups 4 datasets, with a total 46 sequences used for our benchmark. Using ground truth information, object size is divided into 5 categories (0 for *Extremely tiny* up to 4 for *large* [1]) with a mean object size is computed for each sequence (reported on the $x$-axis). To capture displacement between two successive frames, the mean IoU between two successive ground truth bounding boxes is also computed for each sequence (reported on the $y$-axis - BB represents bounding boxes). In contrast to `Falling Object` and `TbD-3D`, both `FMOv2` and `TbD` datasets are more challenging having smaller objects with smaller overlapping successive bounding boxes.

## 2.2 Segment Anything Model 2 (SAM2) Based Trackers

The Efficient Track Anything Model (EfficientTAM) [12], Distractor-Aware Memory for SAM2 (DAM4SAM) [11], and SAM-based Unified and Robust zero-shot

visual tracker with motion-Aware Instance-level memory (SAMURAI) [13] trackers are SAM2 [7] based trackers. A primary reason for choosing SAM2 [7] based trackers is that they do not require training or retraining. Each tracker offers distinct advantages in memory management, performance and adaptation for visual object tracking. Specifically, DAM4SAM is selected to avoid distractors near the target object, particularly in scenarios involving multiple instances of the target object. Meanwhile, SAMURAI is chosen for its integrated motion modeling and motion-aware instance-level memory, which enhance performance in crowded scenes with fast-moving or self-occluding objects, where SAM2 encounters challenges. EfficientTAM is chosen for its ability to deliver accelerated performance and reduced computational costs, making it ideal for efficient video object segmentation across various platforms.

**SAM2.** The image encoder and memory mechanism are essential components of SAM2. The image encoder is responsible for extracting features from frames, while the memory mechanism stores the past $n$ frames to facilitate the segmentation of new frames [12]. This memory mechanism consists of 7 slots for storing 7 frames, with the first slot reserved for the initialized frame. The remaining 6 slots are updated each time a new frame arrives, following a first-in, first-out (FIFO) queue method [12,11,13]. SAM2 generates three output masks and selects the one with the highest predicted Intersection over Union (IoU). However, DAM4SAM [11] noted that simple output masks selection often leads to the inclusion of distractors from previous frames before a tracking failure occurs due to the accumulation of misleading information. Additionally, this straightforward approach can create further issues in crowded scenes where target and background objects have similar appearances. Simply relying on the previous $n$ frames can also result in the storage of misleading features during occlusion [13].

**EfficientTAM.** EfficientTAM offers a lightweight version of SAM2 to reduce the high computational complexity of the image encoder and memory module, particularly for video object segmentation on mobile devices. Unlike the original SAM2, EfficientTAM adopts a lightweight Vision Transformer (ViT) image encoder for improved efficiency. In the memory mechanism, tokens are small pieces of information that the model uses to remember different parts of an image or data. Two adjacent tokens are similar, with a small difference between them, defined by a constant that specifies the acceptable level of similarity. To enhance memory efficiency, EfficientTAM avoids storing multiple nearly identical memory tokens for similar parts of an image or data. Instead, it consolidates these similar tokens into a single representative token. This means that rather than keeping separate tokens for each similar part, the model creates one token that captures the essence of all those similar parts. As a result, the overall set of tokens becomes a coarser representation of the original, meaning it retains the same total number of tokens but simplifies the information they represent. By doing this, EfficientTAM can process information using fewer unique tokens.

This not only speeds up the calculations but also reduces the amount of memory required, making the model more efficient.

**DAM4SAM.** Distractors are elements within the visual field that complicate the tracking of a target object. These can be categorized into two types: external distractors, which are nearby objects that share visual similarities with the target such as an independent instance of the target object, and internal distractors, which are similar regions found on the target itself when only a portion of it is being tracked. The challenge posed by external distractors is particularly pronounced when the target exits and subsequently re-enters the field of view, as these similar-looking objects can lead to confusion in accurately identifying the target. To reduce to distractors failures, the FIFO memory update mechanism of SAM2 has been replaced with a Distractor-Aware Memory (DAM) management strategy. It divides the memory into Recent Appearance Memory (RAM) and Distractor Resolving Memory (DRM), utilizing a new memory management protocol for updates. The 3 slots in the RAM are updated every 5 frames with the FIFO mechanism in case the target object already exists. The DRM, which accounts for the remaining four slots, fixes the first slot in the initialization frame as SAM2.

**SAMURAI.** SAMURAI adapts the SAM2 model to handle distractors and incorporates motion cues for improved memory management. A Kalman filter-based motion modeling is integrated to manage fast-moving and occluded objects in crowded scenes. In addition to the mask affinity score and object occurrence score, the output of the motion modeling, referred to as the motion score, is used to select frames for memory, rather than relying on the n-previous frames as SAM2 does. Instead of relying on a fixed window of frames, a dynamic frame selection process referred to as "Motion-Aware Instance-Level Memory" selectively chooses only the most reliable frames from a sequence to update the memory. A frame is considered a valid candidate for memory if it achieves a good affinity and motion score. If the remaining memory slots have not been filled, the tracked object is considered to be occluded or has disappeared, and frames are filled accordingly. Conversely, frames are discarded if they have a poor affinity or motion score. This approach ensures that the memory is composed of high-confidence frames.

## 3    Method for benchmarking

### 3.1    Pretrained models for benchmarking

All trackers have been initialized via ground-truth bounding boxes and evaluated under their default model configurations: SAM2 and DAM4SAM with the `SAM2.1 Hiera Large (Hiera-L)` model, SAMURAI with the `SAM2.1 Hiera Base Plus (Hiera-B+)` model, and EfficientTAM with the `efficienttam_s` model.

### 3.2    Tracker Initialisation

As observed in previous recent works [1,3], initialization of trackers with bounding boxes performs better than using points to provide an exemplar template to target in the following frames. We have chosen here to initialize all trackers with the first bounding box of the FMOX-labelled object when it occurs in each sequence. Of note, sometimes this first bounding box is not in the first frame of the sequence but occurs in a later frame.

### 3.3    Tracker scoring

The performance metrics Dice and IoU have been chosen here to compare the tracker predicted bounding boxes with the ground truth ones. These metrics are computed on each frame for which FMOX provides a ground truth bounding box. Frames without a ground truth bounding box, for instance when the object has disappeared from view, are not taken into account in the computation of the sequence mean IoU (mIoU) and mean Dice (mDice). IoU, Dice, and their respective means computed over sequence are values between 0 (for object missed or not tracked) and 1 (for perfect detection and tracking).

The first frame used for initialization of the tracker is omitted from performance calculations because its object location is provided by the FMOX ground truth. Furthermore, for frames where the tracker fails to predict a bounding box, the IoU and Dice scores are set to zero, representing the worst possible scores for these frames. These zeros are included as part of the scores computed for each sequence in FMOX (mIoU and mDice).

### 3.4    FMOX for benchmarking trackers

To evaluate the performance of each tracker, we utilise the 46 sequences of the FMOX dataset, none of which were used during the training process. Hence we ensure that no data leakage has occurred between the models and the evaluation dataset.

## 4    Results and discussion

**Quantitative results.** Table 1 provides the minimum, maximum, mean, and median values of mIoU and mDice metrics computed for FMOX dataset for each of the trackers. Both metrics concur in finding the best overall performance with DAM4SAM using both the mean and median (equivalent to a robust mean) computed with all sequences in FMOX. In contrast, EfficientTAM has the worst performance of the four trackers tested.

In Table 2, we present the performance ranking of each tracker across each FMO datasets with box plots presented in Figure 2. Our findings align with those reported in Mosev2 [3], as DAM4SAM consistently outperforms other trackers on the FMO datasets, achieving the highest median and average mIoU and mDice

**Table 1.** Overall results obtained on the 46 sequences in the FMOX dataset. Both the mean and median, computed with the mDICE and mIoU for each sequence in FMOX, show DAM4SAM performing best (in bold font and an asterisk (*)). However the zeros scores observed for the minimum highlight that trackers completely fails for some sequences, while the maximum scores show that sometimes SAM2 outperforms other tracker for some sequences. Values range from [0, 1] (0 = bad, 1 = good).

| mIoU(↑) | SAM2 | EfficientTAM | DAM4SAM | SAMURAI |
|---|---|---|---|---|
| MIN | 0.000 | 0.000 | 0.000 | 0.001 |
| MAX | 0.928 | 0.799 | 0.819 | 0.925 |
| MEDIAN | 0.591 | 0.548 | **0.605*** | 0.596 |
| MEAN | 0.461 | 0.438 | **0.505*** | 0.488 |

| mDice (↑) | SAM2 | EfficientTAM | DAM4SAM | SAMURAI |
|---|---|---|---|---|
| MIN | 0.000 | 0.000 | 0.000 | 0.002 |
| MAX | 0.962 | 0.885 | 0.899 | 0.961 |
| MEDIAN | 0.699 | 0.684 | **0.744*** | 0.736 |
| MEAN | 0.545 | 0.520 | **0.600*** | 0.579 |

scores. This indicates that DAM4SAM delivers both accurate and stable tracking performance across diverse datasets. EfficientTAM ranks lowest overall, with the poorest median and mean scores and frequent missed detections, underscoring its limitations in these challenging scenarios. SAMURAI and SAM2 demonstrates moderate performance, generally outperforming EfficientTAM. The results indicate that a tracker may perform strongly on some subsequences, while others exhibit poor or no performance. As highlighted in [1], initializing trackers with highly motion-blurred frames can adversely affect their performance. Although EfficientTAM has been shown to perform competitively to pipelines dedicated to track fast moving objects [1], its primary design focus is on reducing the computational cost of SAM2, making it more suitable for deployment across various platforms.

**Table 2.** Model performance rankings per dataset in FMOX based on IoU and Dice Score. Only *FMOv2* and *TbD* have sequences with object size extremely tiny to small as per classification provided in FMOX JSON [1]. In addition, ground truth bounding boxes rarely overlap between frames $n$ and $n + 1$ in the *FMOv2* and *TbD* datasets.

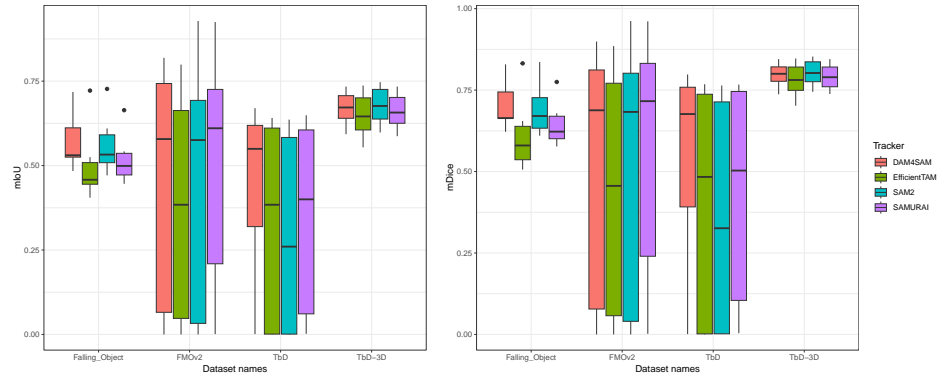| Datasets in FMOX | Ranking (best to worst) |
|---|---|
| *Falling Object* | DAM4SAM > SAM2 > SAMURAI > EfficientTAM |
| *TbD-3D* | SAM2 > DAM4SAM > SAMURAI > EfficientTAM |
| *FMOv2* | SAMURAI > DAM4SAM > SAM2 > EfficientTAM |
| *TbD* | DAM4SAM > SAMURAI > EfficientTAM > SAM2 |

**Fig. 2.** Box plots for mIoU and mDice results on each of the 4 datasets (reported on the x-axis) included in FMOX. Both *FMOv2* and *TbD* are more challenging as objects tracked are of smaller sizes with often non overlapping ground truth bounding boxes between frames $n$ and $n+1$ (cf. Fig. 1). The low minima (=0) highlight the challenge presented by some sequences in these datasets for the trackers tested.

**Compute costs.** The computational workload for each of the four trackers was processed using an NVIDIA GeForce RTX 4090 GPU. All experiments are conducted on a workstation equipped with a 13th Gen Intel Core i9 processor, 64 GB of RAM, CUDA 12.4.1, and Ubuntu 20.04.6 running on Windows Subsystem for Linux (WSL). Table 3 reports computation times: as expected EfficientTAM is the fastest.

**Table 3.** Execution times (in seconds ↓) for the tested trackers across the 4 FMO datasets in FMOX. EfficientTAM offers the lowest computational overhead, ranking as the fastest tracker on all datasets. SAMURAI is the second fastest tracker but also has good accuracy in contrast to EfficientTAM (cf. Tab. 1).

| Tracker / Dataset | DAM4SAM | SAMURAI | EfficientTAM | SAM2 |
|---|---|---|---|---|
| Falling Object | 67.16 | 27.45 | **24.64** | 51.29 |
| FMOv2 | 1316.92 | 515.73 | **410.07** | 890.80 |
| TbD | 261.42 | 93.68 | **67.31** | 394.44 |
| TbD-3D | 168.44 | 56.32 | **39.44** | 203.92 |

**IoU per frames.** To better understand the temporal dynamics of each trackers performance, we analyze the frame-by-frame IoU plots as shown in Figures 3, 4, and 5. These figures highlight the behavior and failure points of each tracker throughout the sequences. The x-axis represents frame numbers derived from the sequence indicated in the frame names. Trackers occasionally fail to detect

objects when they become nearly invisible for one or two frames, most likely due to motion blur. For instance, for the sequence v_rubber_GTgamma in the *Falling Object* dataset, the EfficientTAM tracker exhibited multiple missed detections (see Fig. 3). Notably, it failed to produce a prediction for frames 39, 40, and 41, immediately following a frame with strong motion blur (frame 38). While other trackers show similar transient failures, EfficientTAM is more consistently vulnerable to this prolonged loss of tracking following motion blur events. Similarly, this behavior is observed in Figure 4. While other trackers recovered after just one or two frames of failure, EfficientTAM exhibited a more prolonged failure, missing the object for three consecutive frames (35, 36, and 37). On the other hand, in some sequences, the correlations between the trackers' performances are very strong (see top plot Fig. 5). Conversely, in other sequences, one or a few trackers perform exceptionally well while others perform near zero (see bottom plot Fig. 5).
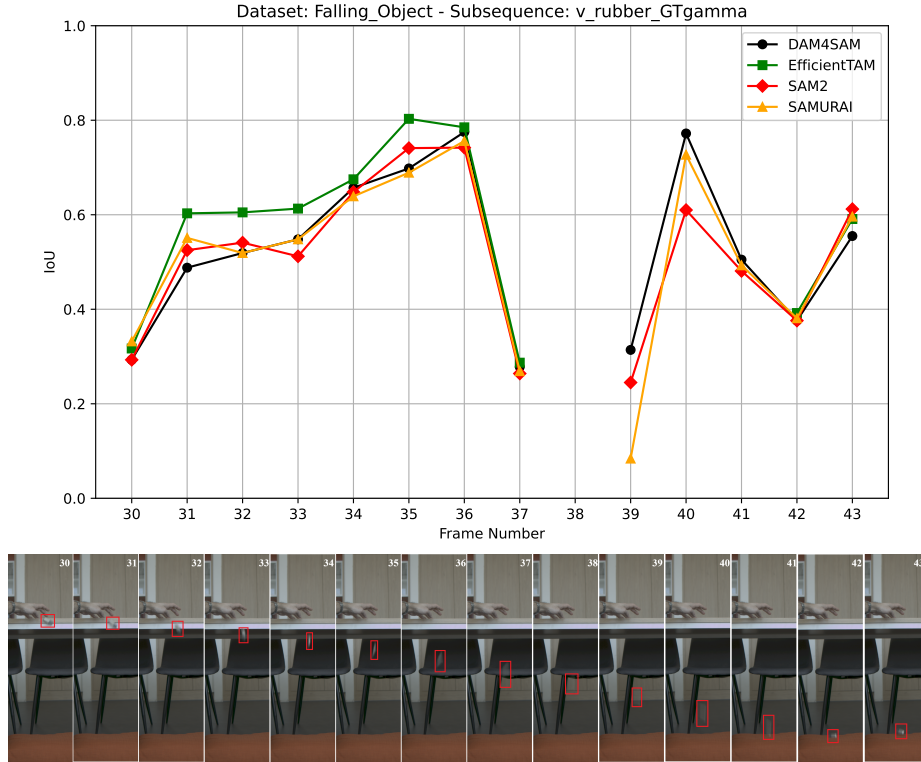


**Fig. 3.** Tracking performance (IoU) across frames on the sequence v_rubber_GTgamma from dataset *Falling Object*. All trackers fails to propose a bounding box for frame 38 while EfficientTAM also fails for frames 39 to 41 included. Corresponding frame numbers are given on top of each frame, and object (rubber) locations are indicated with red ground truth bounding boxes.
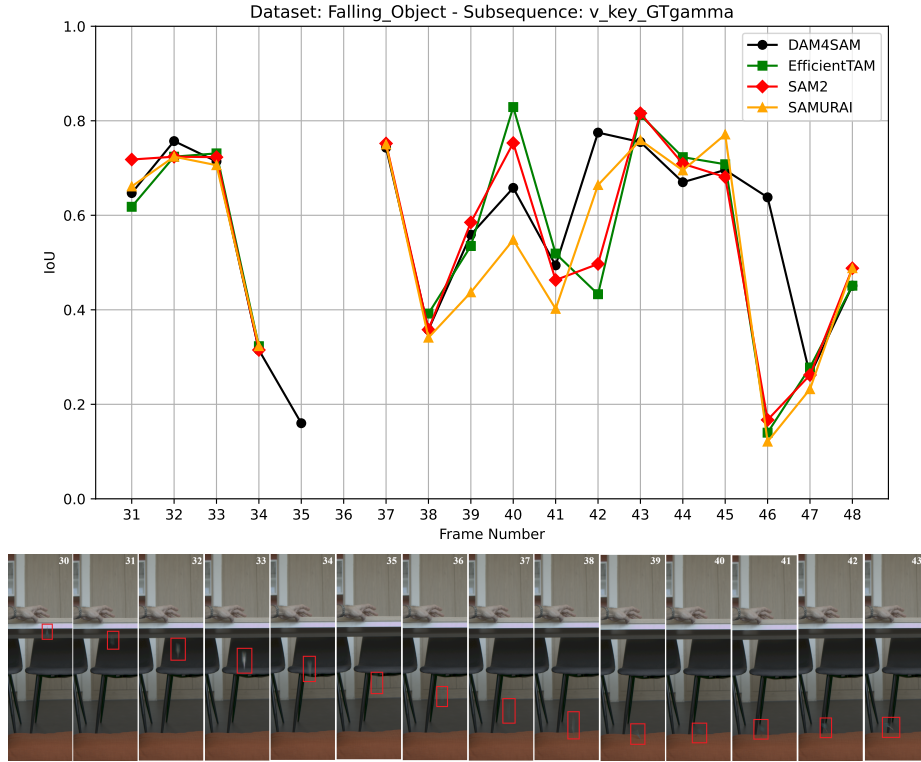
**Fig. 4.** Tracking performance (IoU) across frames on the sequence `v_key_GTgamma` from *Falling Object* dataset. All trackers fail to propose a bounding box for frame 36. For frame 35, DAM4SAM is the sole successful tracker.

## 5   Conclusion

We have benchmarked several trackers on several datasets with fast moving objects, and we have shown that both SAMURAI and DAM4SAM trackers outperform SAM2 and EfficientTAM. Using FMOX classification of object sizes [1] for these datasets, we note that datasets presenting sequences with smaller moving objects (and in addition with non overlapping ground truth bounding boxes between successive frames) affect tracker performance as measured by mIoU and mDice.
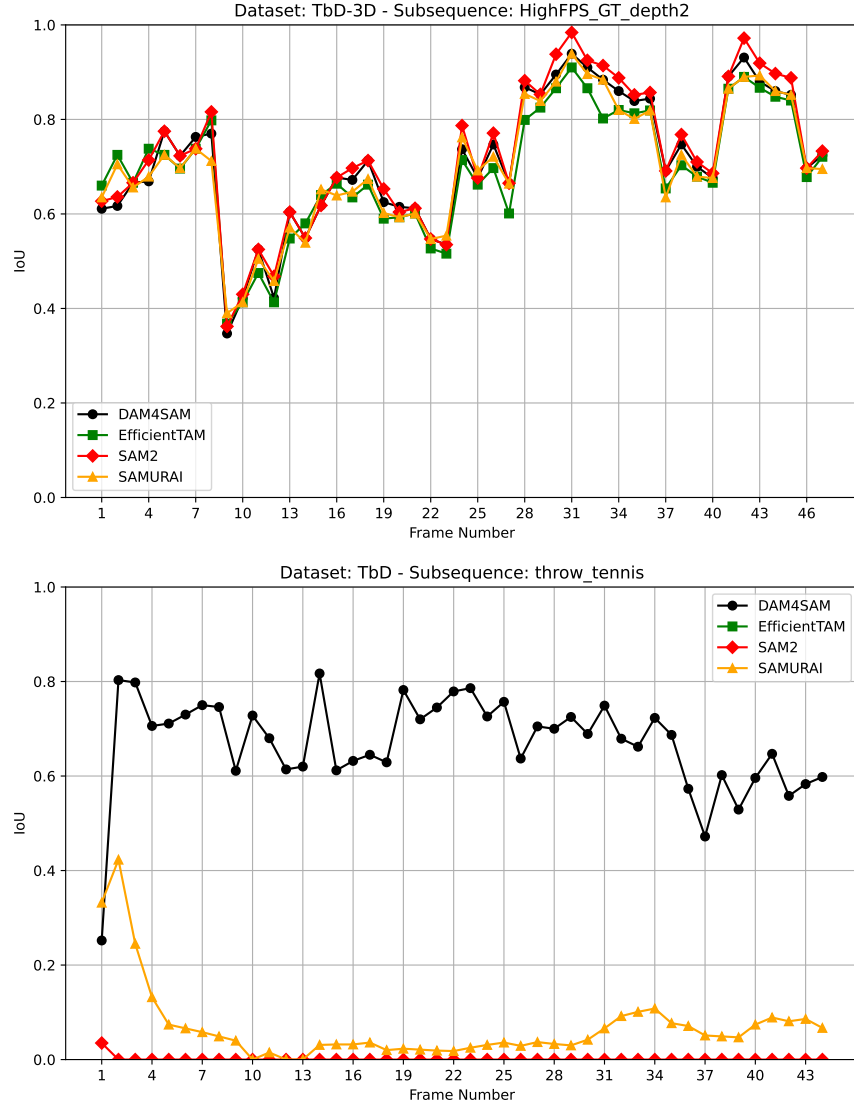
**Fig. 5.** Examples of tracking performance (IoU) across Frames in sequences `HighFPS_GT_depth2` in *TbD-3D* dataset (top: all trackers perform well and provided similar results) and `throw_tennis` from *TbD* dataset (bottom: all trackers performed poorly, with the exception of DAM4SAM; EfficientTAM failed to initialize for tracking due to the strong motion blur present on the object, resulting in no performance curve being generated for this sequence in the graph).

**Disclosure of Interests.** The authors declare they have no competing interests.

# References

1. Aktas, S., Markham, C., McDonald, J., Dahyot, R.: Benchmarking efficienttam on fmo datasets. In: Irish Machine Vision and Image Processing (IMVIP 2025). pp. 59–66. Ulster University, Derry-Londonderry, Northern Ireland (2025). `https://doi.org/10.48550/arXiv.2509.06536`

2. Cui, Y., Jiang, C., Wang, L., Wu, G.: Mixformer: End-to-end tracking with iterative mixed attention. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 13608–13618 (2022)

3. Ding, H., Ying, K., Liu, C., He, S., Jiang, X., Jiang, Y.G., Torr, P.H., Bai, S.: Mosev2: A more challenging dataset for video object segmentation in complex scenes. arXiv preprint arXiv:2508.05630 (2025)

4. Fan, H., Lin, L., Yang, F., Chu, P., Deng, G., Yu, S., Bai, H., Xu, Y., Liao, C., Ling, H.: Lasot: A high-quality benchmark for large-scale single object tracking. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5374–5383 (2019)

5. Kotera, J., Matas, J., Šroubek, F.: Restoration of fast moving objects. IEEE Transactions on Image Processing **29**, 8577–8589 (2020)

6. Kotera, J., Rozumnyi, D., Šroubek, F., Matas, J.: Intra-frame object tracking by deblatting. In: 2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW). pp. 2300–2309 (2019). `https://doi.org/10.1109/ICCVW.2019.00283`

7. Ravi, N., Gabeur, V., Hu, Y.T., Hu, R., Ryali, C., Ma, T., Khedr, H., Rädle, R., Rolland, C., Gustafson, L., Mintun, E., Pan, J., Alwala, K.V., Carion, N., Wu, C.Y., Girshick, R., Dollar, P., Feichtenhofer, C.: SAM 2: Segment anything in images and videos. In: The Thirteenth International Conference on Learning Representations (2025), `https://openreview.net/forum?id=Ha6RTeWMd0`

8. Rozumnyi, D., Kotera, J., Sroubek, F., Matas, J.: Sub-frame appearance and 6d pose estimation of fast moving objects. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6778–6786 (2020)

9. Rozumnyi, D., Kotera, J., Sroubek, F., Novotny, L., Matas, J.: The world of fast moving objects. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5203–5211 (2017)

10. Rozumnyi, D., Kotera, J., Šroubek, F., Matas, J.: Tracking by deblatting. International Journal of Computer Vision **129**(9), 2583–2604 (Jun 2021). `https://doi.org/10.1007/s11263-021-01480-w`

11. Videnovic, J., Lukezic, A., Kristan, M.: A distractor-aware memory for visual object tracking with sam2. In: Proceedings of the Computer Vision and Pattern Recognition Conference. pp. 24255–24264 (2025)

12. Xiong, Y., Zhou, C., Xiang, X., Wu, L., Zhu, C., Liu, Z., Suri, S., Varadarajan, B., Akula, R., Iandola, F., Krishnamoorthi, R., Soran, B., Chandra, V.: Efficient track anything (2024), `https://arxiv.org/abs/2411.18933`, `https://yformer.github.io/efficient-track-anything/`

13. Yang, C.Y., Huang, H.W., Chai, W., Jiang, Z., Hwang, J.N.: Samurai: Adapting segment anything model for zero-shot visual tracking with motion-aware memory. arXiv preprint arXiv:2411.11922 (2024)