

Beyond SOT: Tracking Multiple Generic Objects at Once

Christoph Mayer^{1*} Martin Danelljan¹ Ming-Hsuan Yang² Vittorio Ferrari²
 Luc Van Gool¹ Alina Kuznetsova²
¹ETH Zürich ²Google Research

{christoph.mayer, martin.danelljan, vangool}@vision.ee.ethz.ch {minghsuan, vittoferrari, akuznetsa}@google.com

Abstract

Generic Object Tracking (GOT) is the problem of tracking target objects, specified by bounding boxes in the first frame of a video. While the task has received much attention in the last decades, researchers have almost exclusively focused on the single object setting. However multi-object GOT poses its own challenges and is more attractive in real-world applications. We attribute the lack of research interest into this problem to the absence of suitable benchmarks. In this work, we introduce a new large-scale GOT benchmark, LaGOT, containing multiple annotated target objects per sequence. Our benchmark allows users to tackle key remaining challenges in GOT, aiming to increase robustness and reduce computation through joint tracking of multiple objects simultaneously. In addition, we propose a transformer-based GOT tracker baseline capable of joint processing of multiple objects through shared computation. Our approach achieves a $4\times$ faster run-time in case of 10 concurrent objects compared to tracking each object independently and outperforms existing single object trackers on our new benchmark. In addition, our approach achieves highly competitive results on single-object GOT datasets, setting a new state of the art on TrackingNet with a success rate AUC of 84.4%. Our benchmark, code, results and trained models are available at <https://github.com/visionml/pytracking>.

1. Introduction

Visual object tracking is a fundamental problem in computer vision. Over the years the research effort has been directed mainly to two different task definitions: Generic Object Tracking (GOT) [2, 5, 13, 26, 32, 34, 61] and Multiple Object Tracking (MOT) [6, 16, 23, 56, 69–71]. MOT aims at detecting and tracking all objects from a predefined class category list (see Fig. 1), whereas all other objects are ignored. In contrast, GOT focuses on the scenario where a pri-

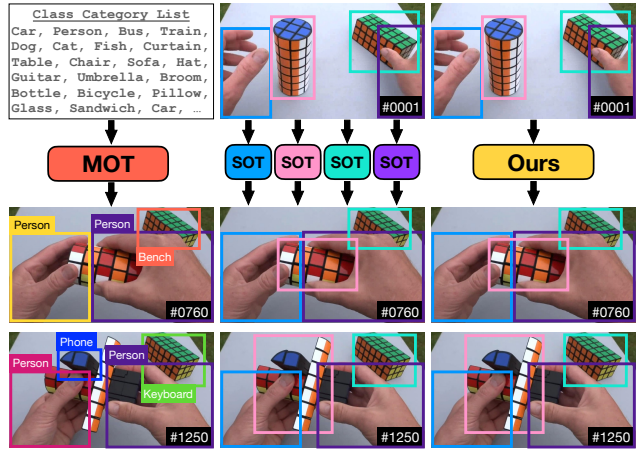


Figure 1. Multiple Object trackers (MOT) track all the objects corresponding to classes in a predefined category list, while all other objects are ignored. Single Object Tracking (SOT) methods focus on tracking only a single user-specified object per video. Thus, when encountered with multiple objects, such methods must resort to independent tracking of each object. This leads to a directly linear increase in computation. Our tracker can track multiple generic objects jointly that are defined via user-specified bounding boxes, leading to the opportunity of computational savings and to exploit inter-object information for improved robustness. The box colors correspond to track IDs.

ori information about the object’s appearance is unknown. Thus, the target model of the object’s appearance must be learned at test time from a single user-specified bounding box in the initial frame, see Fig. 1.

While GOT has a long history of active research, GOT methods and benchmarks focused so far on tracking a single object per video such that the term *Single Object Tracking (SOT)* was introduced. However, the task of GOT is not limited to tracking a single object. In fact, the ability to track multiple generic objects is desired in many real-world applications, such as surveillance, video understanding, semi-automatic video annotation, robotics, and industrial quality control. A method that jointly tracks multiple objects can achieve substantial reduction in computational cost through shared elements, compared to running a separate instance of

*Work done while interning at Google Research.

a SOT method for each object. Moreover, processing multiple targets at the same time has the potential of increasing the robustness of the tracker by joint reasoning.

To facilitate the work on tracking multiple generic objects, we introduce the new multi-object GOT benchmark LaGOT. It provides up to 10 user-specified generic objects in the initial frame visible through the large part of the video. The target objects in one video may correspond to completely different and previously unseen classes. Our benchmark features challenging characteristics such as fast moving objects, frequent occlusions, presence of distractors, camera motion, and camouflage. In total LaGOT contains 528k annotated objects of 102 different classes and an average track length of 71 seconds.

Tracking multiple target objects in the same video poses key challenges and research questions that are typically overlooked by SOT methods. A multi-object GOT method needs to jointly track multiple objects using the first-frame annotations. This could allow the tracker to exploit annotations of potential distractors to improve the robustness of each target model. Furthermore, a joint localization step opens the opportunity for global reasoning across all tracks to reduce the risk of confusing similar objects. Finally, operating on multiple local search area [8, 45, 65] is no longer feasible for a multi-object GOT method because it is inefficient and complicates re-detecting of lost objects.

We tackle these challenges by introducing a new multi object GOT tracker. In order to track all desired target objects at once it operates globally by processing the full frame producing a shared feature representation for all targets. Furthermore, we propose a new generic multiple object encoding that allows us to encode multiple targets within the same training sample. We achieve this by learning a fixed size pool of different object embeddings, each representing a different target. Thus, we query the proposed model predictor with these object embeddings to produce all target models. In addition, we employ a Feature Pyramidal Network (FPN) to increase the overall tracking accuracy while operating on full-frame inputs.

Contributions. (i) We propose a novel large-scale multi-object GOT evaluation benchmark, LaGOT. It provides multiple annotated objects per frame with an average of 2.9 tracks per sequence. We further evaluate several baselines on LaGOT, including two MOT and six SOT methods. We assess their quality by using GOT and MOT metrics.

(ii) We develop a new baseline, TaMOs, a GOT tracker that tracks multiple generic objects at the same time efficiently. To achieve this, we propose a new multi-object encoding, introduce an FPN and apply the tracker globally on the entire video frame. TaMOs demonstrates near constant run-time when increasing the number of targets and operates at an over $4\times$ faster run-time compared to the SOT baselines when tracking 10 objects.

(iii) We analyze TaMOs by assessing the impact of its different components using multiple benchmarks. Furthermore, TaMOs outperforms all baselines on LaGOT, while achieving excellent results on popular SOT benchmarks.

2. Related Work

Object Tracking Benchmarks. Generic object tracking is a well explored topic and many datasets exist. There are specialized datasets and challenges that focus on short-term [21, 28, 30, 32, 48, 61] or long-term tracking [18, 19, 30, 47, 58]. However, all of these benchmarks and datasets share the same setup of only providing a single user-specified bounding box such that only one target is tracked in each video sequence. Recently, GMOT-40 [1] focused on Generic Multi Object Tracking (GMOT), where a single bounding box is provided in the first video frame and all objects that correspond to the same class as the annotated object should be tracked. In contrast to GMOT, we focus on the setting where multiple user-specified targets are given, potentially from different classes.

MOT aims at tracking multiple objects defined by a list of classes and mainly focuses on a single class [16, 29, 56] (usually pedestrians) or on autonomous driving settings, where only a handful of classes are considered [6, 23, 69]. TAO [15] contains objects of a long-tailed class distributions, but provides only sparse annotations due to the costly annotation process. Another related task is open world tracking [39] that aims at detecting and tracking all objects in a video. However, compared to GOT there is no mechanism to guarantee that a specific object is actually detected and tracked. In the Video Object Segmentation (VOS) domain, DAVIS [52] and YouTubeVOS [62] provide multi-object annotations. However, their videos are extremely short (2.9 and 4.5 seconds on average), and are therefore not suitable for tracking. Moreover, the VOS domain provides less challenges for trackers, instead focuses on large objects and a short-term nature, where the predominant challenge is the prediction of accurate fine-grained masks.

Global Generic Object Tracking. Global trackers operate on the whole video frame, rather than in a restricted search area near the object location in the previous frame. This is not only beneficial when tracking multiple objects in the same scene but also facilitates re-detecting lost objects. GlobalTrack [27] and Siam R-CNN [60] track the target by using global RPNs that retrieve target-specific proposals. Recent method for open vocabulary tracking [36] tracks objects of specified classes in MOT fashion by operating on generic RPN proposals. Methods such as MetaUpdater [9] and SPLT [66] operate on local search areas but use a re-detector to re-localize the target if it disappeared from the search area. In contrast, our tracker TaMOs always operates on the entire frame and generates target specific correlation

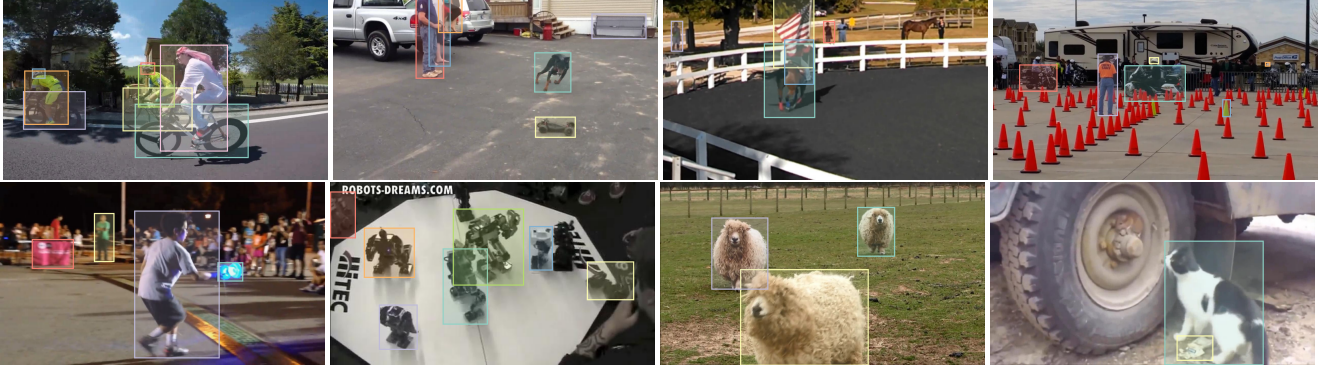


Figure 2. Examples of the annotated objects in the video sequences of our LaGOT dataset. The objects are annotated at 10 FPS. Notice the diversity of the annotated media as well as the complexity of the scenes.

filters instead of target-specific proposals.

Unified Object Tracking. Unified methods aim at tracking both: objects defined by class names or objects defined by a user-specified bounding box. UTT [44] allows to track all pedestrians and one generic object in each video. UTT uses a Transformer to match test frame features with reference features of the detected objects in the initial or previous frame. Unicorn [64] allows to either perform the SOT or the MOT task with the same model and weights solely by varying the input data type. In contrast, our method tracks multiple generic objects at the same time instead of one generic object or multiple objects of known classes.

Transformers for Generic Object Tracking. Tracking has seen a tremendous progress in recent years with the advent of Transformers [59]. Most such trackers share the idea of fusing the search area and the template image features by using a Transformer [7, 8, 45, 65, 67, 68]. MixFormer [8] and OTrack [67] employ a Transformer to jointly extract and fuse the template and search area features. TransT [7], STARK [65], SwinTrack [37] and ToMP [45] use a backbone to extract features and employ cross attention to fuse the feature representations. However, none of these trackers can easily be extended to jointly track multiple objects, which is addressed in this work.

3. LaGOT Benchmark

In this section we first introduce the multi-object GOT task and discuss its differences to other object tracking tasks. Then, we introduce our new benchmark LaGOT.

3.1. Multi-object GOT Task

Multi-object GOT is the task of tracking multiple generic target objects in a video sequence. The target objects are defined by user-specified bounding boxes in the initial frame of the video. Thus, the target objects are generic in the sense that their class category is unknown and there might be no object of the same category in the training data, see Fig. 1.

Multi-object GOT vs. SOT. SOT requires to track only

a single target object defined by the user [19, 32, 61], whereas multi-object GOT focuses on tracking multiple user-specified generic target objects in the same video.

Multi-object GOT vs. MOT. (i) The MOT task requires to track all objects of known classes, whereas for multi-object GOT target objects in each video are defined by user-specified boxes. Consequently, multi-object GOT is a one-shot problem where the target objects are unknown at training time and are only available during inference. In contrast, traditionally MOT methods track all objects corresponding to the categories defined at training time. (ii) For the multi-object GOT task an object-id switch is equivalent to a complete failure since the user-specified object is no longer recoverable [43, 61]. Conversely, for MOT methods object-id switches are considered less problematic and are penalized less drastically by the MOT metrics [42].

Multi-object GOT vs. GMOT. GMOT focuses on tracking multiple objects of a single generic object class in each video. The class is defined by a single user-specified bounding box in the initial video frame [1, 20]. Thus, in contrast to multi-object GOT, a GMOT method is unable to track multiple objects of different categories in the same video.

3.2. LaGOT

Benchmark Construction. LaSOT [19] contains diverse and relatively long videos (2430 frames or 81 seconds on average) with challenging tracking scenarios including fast moving objects, camera motion, various object sizes, frequent object occlusions, scale changes, motion blur, camouflage and objects that go out of view or change their appearance. LaSOT provides annotations for a single object in each video but typically multiple objects are present throughout the full sequence and are fairly difficult to track, which is desirable for long-term tracking scenarios. Thus, instead of collecting new videos, we used the popular LaSOT evaluation set and add new annotations for multiple objects in each sequence.

Another large-scale video dataset we considered is

Table 1. Comparison of LaGOT with existing benchmarks that focus on related tasks to multi-object GOT.

Dataset	Task	Object Definition	Num Classes per Video	Tracking Metrics	Num Classes	Num Videos	Avg Video Length (s)	Avg Track Length (num anno.)	Avg Tracks per Video	Num Annotations	Annotation Frequency
TAO val [15]	MOT	class list	≥ 1	Track-mAP	302	988	33.5	21	5.55	115k	1 FPS
GMOT-40 [1]	GMOT	1 box	1	MOTA/IDF1	30	40	10	133	50.65	486k	24-30 FPS
LaSOT val [19]	SOT	1 box	1	Success AUC	70	280	81	2430	1	680k	30 FPS
LaGOT	GOT	n boxes	≥ 1	F1-Score	102	294	75.3	707	2.89	528k	10 FPS

TAO [15] and GMOT-40 [1]. However, compared to LaSOT, TAO contains shorter videos with an average of 33 seconds and its outdoor and road sequences mainly focus on pedestrians and vehicles (60% of all objects in TAO). While the indoor sequences contain rarer object categories, they are often static and are only visible for a short time. Furthermore, TAO contains only sparse annotations (1 FPS). For all these reasons, we used LaSOT instead of TAO to build our benchmark. GMOT-40 [1] contains dense annotations, but videos often contain many objects of a single class. Furthermore, GMOT-40 consists of only 40 short sequences (avg 240 frames or 8 seconds) rendering only 10 different object classes, see Tab. 1. Thus, GMOT-40 is unsuitable to serve as a multi-object GOT benchmark.

Annotation Protocol. First, we inspect all 280 sequences in LaSOT and identify in each video challenging target objects that play an active role and meet the previously specified criteria. Next, we entrust professional annotators to annotate the selected objects in all sequences on every third frame, leading to an annotation frequency of 10 FPS. They use an interactive annotation tool which incorporates an object tracker to speed up the annotation process [33]. A group of researchers verifies the newly obtained annotations and sends low-quality annotations back for correction until all annotations meet our high quality standards. Finally, we post-process the annotations to construct the final tracks. First, we remove all tracks shorter than 4 seconds. Second, we define the starting frame by manually selecting the earliest frame where as many annotated objects as possible are clearly visible. Third, it is not always possible to unambiguously associate all object identities over time due to occlusions and out-of-view events — hence, we either remove ambiguous annotations or cut these videos into multiple sub-sequences, where the object association is clear. We follow this protocol to guarantee a high annotation quality, see Fig. 2 for annotated example frames.

Statistics. Our benchmark LaGOT has 294 videos with 850 tracks leading to over 528k annotated objects. Thus, we almost triple the number of tracks compared to the original LaSOT validation set (and the corresponding evaluation time from 378 to 1006 min). Furthermore, we add 31 additional generic object classes, *e.g.* propeller, tires or fabric bag. We compare the proposed benchmark with the most closely related benchmarks in Tab. 1 (and with many more Tab. 2 in suppl. material). Overall our benchmark con-

tains $10\times$ more class categories than GMOT-40. The average track length of LaGOT is 2121 frames (707 annotated frames), which is $3\times$ longer than in TAO, and almost $10\times$ longer than in GMOT-40.

Annotation Frequency. According to Valmadre *et al.* [58] it is more effective to spend a fixed annotation budget on many videos with sparse box annotations than on fewer videos with dense labels. Thus, we annotate every third frame to reduce the overall annotation cost. To analyze the difference between 10 and 30 FPS annotations, we evaluate five recent trackers on the tracks borrowed from LaSOT, where 30 FPS annotations are available. The mean relative error of the success rate AUC is only 0.237%. This shows that 10 FPS is sufficient on large-scale datasets such as LaSOT and LaGOT, leading to only minor score deviations.

4. Method

In this section we present our tracker TaMOs, which employs a Transformer to jointly model and track a set of arbitrary objects defined in the initial frame of a video. We start from ToMP [45], a recent Transformer-based generic single object tracker that operates on local search area cropped from the full frame, as almost all SOT trackers. ToMP employs a transformer to predict a correlation filter (target model) from the target appearance in the initial frame conditioned on the new frame; the predicted target models is later used to localize the target in the subsequent frames. In Sec. 4.1 we introduce the proposed Transformer-based multi-object tracking architecture and in Sec. 4.2 we discuss the used training protocol.

4.1. Generic Multi-Object Tracker - Overview

An overview of the proposed generic multi-object tracker TaMOs is presented in Fig. 3. First, unlike original ToMP, our tracker operates on the full train and test images instead of crops. The target object encoder uses a pool of learnable object embeddings to encode the location and extent of each target object within a single shared feature map (Sec. 4.1.1). The randomly sampled object embedding then represents a particular target in the entire video sequence: we use the object embedding to condition the model predictor to produce the target model that localizes the target object in the test frame (Sec. 4.1.2). Since operating on the entire video frame increases the computational cost of the Transformer operations, we are limited to a certain fea-

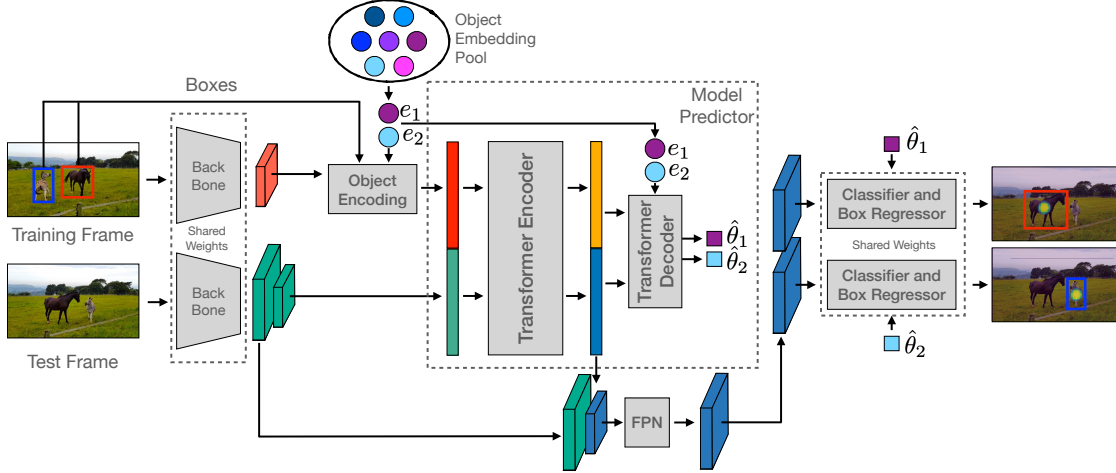


Figure 3. Overview of our tracker TaMOs for joint tracking of multiple targets. First, we extract features from training and test frames. All objects in the training frame are encoded jointly with a multi-object encoding and passed to the model predictor together with the training frame features. The model predictor produces target models $\hat{\theta}_i$ together with enhanced test features. We apply an FPN on the enhanced output features to generate higher resolution test features. Finally, we predict the bounding box of each target by applying the target model $\hat{\theta}_i$ for each target.

ture resolution. To track small objects we propose an FPN-based feature fusion of the test frame features produced by the Transformer with the higher resolution backbone features. We adopt the correlation filter based target localization and bounding box regression mechanism of ToMP but apply both on the higher resolution FPN features instead of the output features of the Transformer (Sec. 4.1.3).

4.1.1 Generic Multiple Object Encoding

To track several target objects efficiently, we propose a novel object encoding to embed multiple objects in a shared feature map without requiring multiple templates. In particular, we extend the single object encoding formulation of ToMP to be applicable for multiple objects. The idea is to replace the foreground embedding with multiple object embeddings, each representing a different target object. Thus, we create a pool $E \in \mathbb{R}^{m \times c}$ of m object embeddings $e_i \in \mathbb{R}^{1 \times c}$. Then, we sample for each target object a random object embedding from the pool E without replacement. Next, we combine the object embeddings with the Gaussian score map $y_i \in \mathbb{R}^{h \times w \times 1}$ that represents the center location of the target object i and the LTRB [57, 63] bounding box encoding $b_i^{\text{ltrb}} \in \mathbb{R}^{h \times w \times 4}$. The final encoding is thus:

$$f_{\text{train}}^{\text{enc}} = f_{\text{train}} + \sum_{i=0}^n e_i \cdot y_i + \sum_{i=0}^n e_i \cdot \phi(b_i^{\text{ltrb}}), \quad (1)$$

where $f_{\text{train}} \in \mathbb{R}^{h \times w \times c}$ are visual features extracted from the full training frame, ϕ is a Multi-Layer Perceptron (MLP) and $n \leq m$ is the number of tracked objects. Note, that in contrast to the object encoding in ToMP, we not only use the object embedding to encode the Gaussian score map

but also the bounding box representation. The object embeddings e_i are learned during training such that the model is able to disentangle the shared feature representation and can identify each object in the training and test features. Note, that the products in Eq. (1) employ multiplications with broadcasting across every dimension whereas the latter uses channel-wise multiplication with broadcasting across the spatial dimensions.

4.1.2 Joint Model Prediction

Now that the target object locations and extents are embedded in the training features, we require a model predictor to produce a target model for each encoded object. The target models are then used to localize the targets in the test frame and to regress their bounding boxes. In order to easily associate the different targets over time, we require a model predictor that can be conditioned on the targets encoded through object embeddings e_i . Furthermore, the model needs to be able to produce all target models jointly to increase the efficiency.

We extend the single target model predictor of ToMP by keeping the Transformer encoder unchanged but by modifying the Transformer decoder. In particular, we query the Transformer decoder with multiple object embeddings e_i at the same time instead of a single foreground embedding,

$$[\hat{\theta}_1, \dots, \hat{\theta}_n] = T_{\text{dec}}([h_{\text{train}}, h_{\text{test}}], [e_1, \dots, e_n]). \quad (2)$$

Here, $\hat{\theta}_i \in \mathbb{R}^c$ is the target model, n is the number of target objects encoded in the training frame and $h_{\text{train}}, h_{\text{test}}$ are the refined output features of the Transformer encoder for the train and test frame.

4.1.3 Target Localization and Box Regression

We use the generated target models to localize the targets and to regress their bounding boxes. We produce a correlation filter for target classification and adopt the bounding box regression branch of ToMP [45]. But instead of applying the target classifier and box regressor on the low-resolution test features h_{test} of the Transformer encoder, we use high resolution features generated with an FPN $\psi(\cdot)$ and obtain the high-resolution multi-channel score map:

$$\hat{y}_i^{\text{high}} = w_i^{\text{cls}}(\hat{\theta}_i) * \psi(h_{\text{test}}, f_{\text{test}}^{\text{high}}), \quad 0 \leq i < n, \quad (3)$$

where $f_{\text{test}}^{\text{high}} \in \mathbb{R}^{2h \times 2w \times c}$ are the high-resolution test features extracted at an earlier stage of the backbone, $w_i^{\text{cls}}(\hat{\theta}_i)$ refers to the discriminative correlation filter for the target object i obtained from the predicted target model $\hat{\theta}_i$. Similarly we obtain the high-resolution multi-channel bounding box regression maps b_i^{high} .

4.2. Training

During training we employ a classification and a bounding box regression loss. We compute both losses for the predictions obtained by processing each FPN feature map (low-res and high-res) as well as the output test features h_{test} of the Transformer encoder. The classification loss is

$$L_{\text{cls}} = \sum_{i=0}^n L_{\text{focal}}(\hat{y}_i, y) + \sum_{j=n}^m L_{\text{focal}}(\hat{y}_j, 0), \quad (4)$$

Here we assume that the first n object embeddings e_i were used to encode the n objects marked in the training frame whereas the remaining $m - n$ object embeddings were not used to encode any objects. Thus, we require that the resulting score maps \hat{y}_j that correspond to an unused object embedding e_j produce low scores everywhere (second sum in Eq. (4)). This step tightly couples the object encoding and decoding. Omitting this term not only decreases the overall performance but slows down the training progress.

In contrast to classification, we enforce the generalized IoU-Loss [54] for bounding box regression only for the predictions that actually correspond to an encoded object and ignore those corresponding to unused object embeddings.

Training Details. We randomly sample an image pair consisting of one training and one test frame from a training video. The frames are re-scaled and padded to a resolution of 384×576 . We train our tracker on the training splits of LaSOT [19], GOT10k [28], TrackingNet [48], MS-COCO [38], ImageNet-Vid [55], TAO [15], and YoutubeVOS [62]. Note, that we remove all videos from the TAO training set that overlap with the evaluation set of LaSOT. We randomly sample for each epoch 40k image pairs with equal probability from all datasets. In order to leverage SOT datasets and training all object embeddings e_i equally,

we assign random object ids to all objects in the sampled training pair. Note, that both SOT and MOT datasets are crucial to train the proposed tracker. Without MOT datasets the tracker is unable to learn multiple target models at the same time and avoiding SOT datasets leads to inferior tracking quality. We train the tracker for 300 epochs on 4 Nvidia A100 GPUs. Our method is implemented using PyTracking [12] (see suppl. material for further details).

5. Experiments

To illustrate the challenges of our proposed GOT benchmark, we evaluate several recent trackers along with our proposed tracker TaMOs on LaGOT (Sec. 5.1). In addition, we compare TaMOs to recent trackers on several SOT benchmarks (Sec. 5.2) and present an ablation study (Sec. 5.3), evaluating the impact of different components of our tracker.

5.1. State-of-the-Art Evaluation on LaGOT

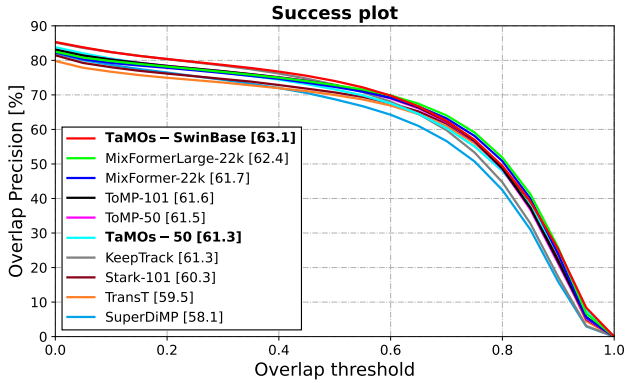
We evaluate our tracker with a ResNet-50 and a Swin-Base backbone as well as six single object trackers (SuperDiMP [12], KeepTrack [46], TransT [7], STARK [65], ToMP [45], and MixFormer [8]) and two multi object trackers (QDTrack [50] and OVTrack [36]) on LaGOT.

Metrics. We measure the performance of a tracker in the One Pass Evaluation (OPE) setting. The standard GOT Success rate Area Under the Curve (AUC) metric [18, 19, 21, 47–49, 61]. does not account for false positive predictions when a target gets occluded or is out of view. While this is not a big issue in standard SOT datasets, where the target object is present in the vast majority of frames, it becomes vital in long-term tracking. In LaGOT objects are more frequently invisible due to occlusions or moving out-of-view. To capture this aspect, we employ the VOTLT [31, 43] metric that penalizes false positives. It computes the IoU-weighted precision-recall curve and ranks the trackers according to their F1-score.

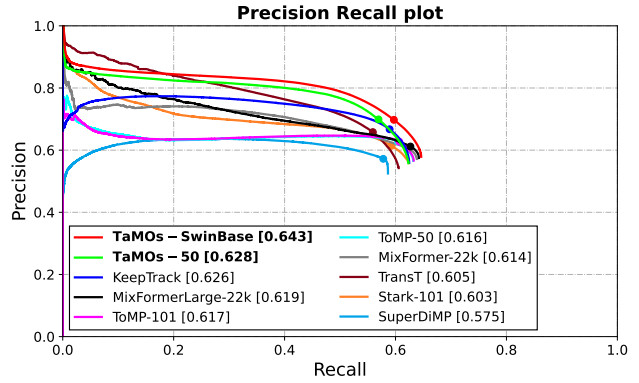
5.1.1 Comparison to SOT Methods

SOT trackers are limited to track only a single target at once. Thus, multiple instances of the same tracker need to be run in parallel to track multiple objects in the same sequence leading to a linearly increasing run-time, see Fig. 1.

Results. Fig. 4a shows the success rate of all trackers on LaGOT. We observe that SOT trackers perform well on LaGOT. However, our multi-object tracker TaMOs achieves the best AUC, even outperforming the state-of-the-art SOT tracker MixFormerLarge-22k [8]. We further observe that TaMOs is as robust as KeepTrack [46] ($T < 0.4$), where the gap to the remaining trackers is particularly prominent. This demonstrates the potential of a global multiple object GOT method. Fig. 4b shows the *tracking* Precision-Recall curve



(a) Success Plot



(b) Precision-Recall Plot

Figure 4. Success plot, showing OP_T , on LaGOT (AUC is reported in the legend). Tracking Precision-Recall curve on LaGOT – VOTLT is reported in the legend (the highest F1-score).

on LaGOT. Both versions of TaMOs outperform all other SOT trackers. The highly robust object presence scores predicted by our tracker lead to a superior precision at all recall rates > 0.2 . Moreover, our approach achieves the best maximal recall and outperforms all previous methods in VOTLT by 1.7 points. This demonstrates that joint tracking of multiple objects and global search benefit the object localization and identification capabilities of the tracker. For further insights we show MOT metrics on LaGOT in Tab. 3. Our tracker achieves the best results for every MOT metric and outperforms MixFormerLarge-22k by 5.9 points in MOTA.

Run-Time Analysis. We evaluate the run-time on a single A100 GPU. Tab. 2 reports a run-time analysis of our tracker TaMOs compared to ToMP, with both employing a ResNet-50 backbone. While TaMOs is slower than ToMP for a single object, due to the higher resolution required for full-frame tracking, our approach already reaches an advantage for 2 concurrent objects. As ToMP needs to run a separate independent tracker for each new object, our approach achieves a $4\times$ speedup for 10 concurrent objects. Furthermore, the analysis demonstrates that TaMOs achieves almost a constant run-time even when increasing the number of targets. TaMOs-SwinBase achieves 13.1 FPS for a single object and 9.3 FPS when jointly tracking 10 objects.

5.1.2 Comparison to MOT Methods

MOT methods are designed to track multiple objects in a video sequence and are thus used as baselines for LaGOT. However, in MOT the targets are defined via a list of classes whereas in multi-object GOT targets are defined by user-specified bounding boxes in the initial frame. Hence, to be able to track generic objects MOT methods need to be trained on large vocabulary datasets — then we can greedily match the detected tracks with the bounding boxes on the initial frame to track user-specified objects. Alternatively, the recent open-vocabulary MOT method OVTrack [36] allows to track objects of any class. We employ QDTrack [50]

Table 2. Run-time analysis (in FPS) between our baseline model ToMP and our tracker TaMOs.

	1 Object	2 Objects	5 Objects	10 Objects
ToMP-50	34.7	17.4	7.0	3.4
TaMOs-50	19.2	17.9	16.3	13.9

Table 3. Comparison of GOT and MOT metrics on LaGOT.

		F1-Score	Success	HOTA	MOTA	IDF1	OWTA
GOT	TaMOs-SwinBase	0.643	63.1	62.1	58.2	74.7	68.9
	TaMOs-50	0.628	61.3	60.0	52.9	72.0	67.1
SOT	MixFormerLarge-22k	0.619	62.4	61.5	52.3	74.3	69.0
	ToMP-101	0.617	61.6	60.1	51.9	73.8	67.5
	STARK-101	0.603	60.3	59.4	49.0	72.5	67.0
	TransT	0.605	59.5	57.7	46.6	70.7	65.6
	KeepTrack	0.626	61.3	59.1	51.3	73.8	66.2
	SuperDiMP	0.575	58.1	56.1	43.2	69.7	63.8
MOT	QDTrack	0.187	19.2	22.2	-115.8	16.3	36.3
	OVTrack	0.128	13.4	24.4	13.9	23.5	25.9

and open-vocabulary OVTrack [36] as MOT baselines. QDTrack is trained on LVIS [24] and TAO. We provide OVTrack in each video with the class name of the target.

Results. QDTrack and OVTrack achieve a VOTLT F1-Score of 0.187 and 0.128 respectively, performing inferior to all other trackers. Neither of the MOT trackers is robust enough and both fail to track rare or unknown generic objects. To further explore the limitations of MOT methods in our setting, we evaluate ‘Oracle’ versions, where we select the track ID that maximizes the scores on LaGOT. Even with such oracle information, the performance of QDTrack and OVTrack is by far inferior to any evaluated SOT baseline (VOTLT 33.1 and 23.0 respectively). In addition we evaluate both trackers using all its predicted tracks with MOT metrics, see Tab. 3. QDTrack tracks multiple background objects that are not annotated in LaGOT leading to many False Positives (FPs), and OVTrack tracks unannotated objects as well since the videos are not annotated exhaustively on class levels. Thus, traditional MOT tracking metrics such as MOTA, HOTA and IDF1 are unsuitable to evaluate MOT trackers on LaGOT. Instead, we concentrate

on the OWTA metric [39] that focuses on Detection Recall and Association Accuracy and thus ignores FPs. QDTrack achieves 36.3 and OVTrack 25.9, which are still the lowest OWTA scores compared to SOT and GOT trackers.

5.2. State-of-the-Art Comparison on SOT Datasets

While TaMOs is built to track multiple objects in a video it can as well track only a single generic object. Thus, we evaluate TaMOs on popular large-scale SOT benchmarks. We deploy the very same tracker in these settings, without altering its weights or any hyper-parameters.

LaSOT [19]. This large-scale dataset consists of 280 test sequences with 2500 frames on average. Tab. 4 shows a comparison to recent SOT trackers. While primarily designed to cope with multiple objects, our tracker achieves the highest precision and the third highest success rate AUC. Note, that neither MixFormer, SwinTrack nor OS-Track operate on the entire video frame, but rely on a local search area to produce such high tracking accuracy.

TrackingNet [48]. This dataset consists of 511 test sequences and predictions are evaluated on a server. Tab. 4 shows that our tracker with SwinBase sets the new state of the art on TrackingNet in terms of success rate and precision AUC. Similarly, our tracker with ResNet-50 achieves the best results among all trackers using that backbone.

The results on both benchmarks show the great potential of applying trackers *globally* without motion priors, such as search area selection [4, 8, 67] or spatial windowing [34, 35].

5.3. Ablation Study

The ablation experiments shown in Tabs. 5 and 6 are performed before the final annotation verification step such that the results compared to the numbers above slightly differ.

Generic Multiple Object Encoding. Tab. 5 shows the effect of the Gaussian score map encoding, the LTRB bounding box encoding and the total number of object embeddings m stored in the pool E . The first two rows in Tab. 5 show that the LTRB encoding is more important than the Gaussian encoding (as removing LTRB decreases all results more significantly). Another key factor is the number of different object embeddings, that sets an upper limit on the number of objects that can be tracked. LaGOT requires at least 10 embeddings and our tracker achieves the best results when using a pool size of 10. Increasing the number of embeddings decreases the overall tracking performance.

Architecture. Tab. 6 shows that using SwinBase increases the tracking performance on LaSOT and LaGOT. Similarly, adding an FPN improves the results.

Inference. During inference we update the memory by adding a second dynamic training frame similar to ToMP [45]. Since the ground truth bounding boxes are not available, we use the predicted boxes as annotations. We replace the dynamic training frame (update the memory) if the

Table 4. State-of-the-art comparison on SOT datasets.

Method	Venue	LaSOT [19]			TrackingNet [48]		
		Prec	N-Prec	Succ	Prec	N-Prec	Succ
TaMOs-SwinBase	WACV'24	77.8	79.3	70.2	84.2	88.7	84.4
TaMOs-50	WACV'24	75.0	77.2	67.9	82.0	87.2	82.7
SwinTrack [37]	NIPS'22	76.5	—	71.3	82.0	—	84.0
Unicorn [64]	ECCV'22	74.1	76.6	68.5	82.2	86.4	83.0
AiATrack [22]	ECCV'22	73.8	79.4	69.0	80.4	87.8	82.7
OSTrack [67]	ECCV'22	77.6	81.1	71.1	83.2	88.5	83.9
RTS [51]	ECCV'22	73.7	76.2	69.7	79.4	86.0	81.6
MixFormer [8]	CVPR'22	76.3	79.9	70.1	83.1	88.9	83.9
ToMP [45]	CVPR'22	73.5	79.2	68.5	78.9	86.4	81.5
UTT [44]	CVPR'22	67.2	—	64.6	77.0	—	79.7
KeepTrack [46]	ICCV'21	70.2	77.2	67.1	73.8	83.5	78.1
STARK [65]	ICCV'21	72.2	77.0	67.1	—	86.9	82.0
TransT [7]	CVPR'21	69.0	73.8	64.9	80.3	86.7	81.4
SuperDiMP [14]	CVPR'20	65.3	72.2	63.1	73.3	83.5	78.1

Table 5. Analysis of different object encoding settings. All tested configurations are not employing the FPN.

Gaussian Encoding	LTRB Encoding	Object Embedding Pool size m	LaSOT AUC	LaGOT AUC	F1
✓	✗	10	58.3	54.0	0.552
✗	✓	10	66.3	60.2	0.620
✓	✓	10	67.2	61.6	0.633
✓	✓	15	65.7	60.0	0.617
✓	✓	20	65.7	58.9	0.603
✓	✓	50	63.1	57.4	0.587

Table 6. Architecture and memory update analysis.

Backbone	FPN	Memory Update	LaSOT AUC	LaGOT AUC	F1
Resnet-50	✗	✓	67.2	60.4	0.621
Resnet-50	✓	✗	66.0	60.2	0.620
Resnet-50	✓	✓	67.9	61.6	0.633
SwinBase	✗	✓	69.5	62.4	0.643
SwinBase	✓	✗	67.9	62.1	0.636
SwinBase	✓	✓	70.2	63.5	0.649

maximal value in each target score map is above the threshold of $\tau = 0.85$. The results in Tab. 6 show that adding a second training frame improves the results on both datasets.

6. Conclusion

We propose a novel multiple object GOT tracking benchmark, LaGOT, that allows to evaluate GOT methods that can jointly track multiple targets in the same sequence. We demonstrate that the proposed task and benchmark are challenging for existing SOT and MOT trackers. We further propose a Transformer-based tracker capable of processing multiple targets at the same time, with a novel generic multi object encoding and an FPN in order to achieve full frame tracking. Our method outperforms recent trackers on the LaGOT benchmark, while operating $4\times$ faster than the SOT baseline when tracking 10 objects. Lastly, our approach also achieves excellent results on popular SOT benchmarks.

Funding: This work was done at Google Research.

Acknowledgement: The authors thank Paul Voigtlaender for his support, fruitful discussions and valuable feedback.

References

- [1] Hexin Bai, Wensheng Cheng, Peng Chu, Juehuan Liu, Kai Zhang, and Haibin Ling. Gmot-40: A benchmark for generic multiple object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6719–6728, June 2021. [2](#), [3](#), [4](#), [18](#)
- [2] Luca Bertinetto, Jack Valmadre, João F Henriques, Andrea Vedaldi, and Philip HS Torr. Fully-convolutional siamese networks for object tracking. In *Proceedings of the European Conference on Computer Vision Workshops (ECCVW)*, October 2016. [1](#)
- [3] Alex Bewley, Zongyuan Ge, Lionel Ott, Fabio Ramos, and Ben Upcroft. Simple online and realtime tracking. In *2016 IEEE international conference on image processing (ICIP)*, pages 3464–3468. IEEE, 2016. [13](#), [14](#)
- [4] Goutam Bhat, Martin Danelljan, Luc Van Gool, and Radu Timofte. Learning discriminative model prediction for tracking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. [8](#), [15](#)
- [5] David S. Bolme, J. Ross Beveridge, Bruce A. Draper, and Yui Man Lui. Visual object tracking using adaptive correlation filters. In *CVPR*, 2010. [1](#)
- [6] Ming-Fang Chang, John Lambert, Patsorn Sangkloy, Jagjeet Singh, Slawomir Bak, Andrew Hartnett, De Wang, Peter Carr, Simon Lucey, Deva Ramanan, and James Hays. Argoverse: 3d tracking and forecasting with rich maps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [1](#), [2](#)
- [7] Xin Chen, Bin Yan, Jiawen Zhu, Dong Wang, Xiaoyun Yang, and Huchuan Lu. Transformer tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021. [3](#), [6](#), [8](#)
- [8] Yutao Cui, Cheng Jiang, Limin Wang, and Gangshan Wu. Mixformer: End-to-end tracking with iterative mixed attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13608–13618, June 2022. [2](#), [3](#), [6](#), [8](#)
- [9] Kenan Dai, Yunhua Zhang, Dong Wang, Jianhua Li, Huchuan Lu, and Xiaoyun Yang. High-performance long-term tracking with meta-updater. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [2](#)
- [10] Kenan Dai, Yunhua Zhang, Dong Wang, Jianhua Li, Huchuan Lu, and Xiaoyun Yang. High-performance long-term tracking with meta-updater. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [15](#)
- [11] Martin Danelljan, Goutam Bhat, Fahad Shahbaz Khan, and Michael Felsberg. ATOM: Accurate tracking by overlap maximization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [15](#)
- [12] Martin Danelljan, Goutam Bhat, and Christoph Mayer. PyTracking: Visual tracking library based on PyTorch. <https://github.com/visionml/pytracking>, 2019. Accessed: 1/07/2022. [6](#)
- [13] Martin Danelljan, Goutam Bhat, Fahad Shahbaz Khan, and Michael Felsberg. ECO: efficient convolution operators for tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2017. [1](#)
- [14] Martin Danelljan, Luc Van Gool, and Radu Timofte. Probabilistic regression for visual tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [8](#), [15](#)
- [15] Achal Dave, Tarasha Khurana, Pavel Tokmakov, Cordelia Schmid, and Deva Ramanan. Tao: A large-scale benchmark for tracking any object. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 436–454. Springer International Publishing, 2020. [2](#), [4](#), [6](#), [18](#)
- [16] Patrick Dendorfer, Aljosa Osep, Anton Milan, Konrad Schindler, Daniel Cremers, Ian Reid, Stefan Roth, and Laura Leal-Taixé. Motchallenge: A benchmark for single-camera multiple target tracking. *International Journal of Computer Vision (IJCV)*, 129(4):1–37, 2020. [1](#), [2](#), [18](#)
- [17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009. [18](#)
- [18] Heng Fan, Hexin Bai, Liting Lin, Fan Yang, Peng Chu, Ge Deng, Sijia Yu, Mingzhen Huang, Juehuan Liu, Yong Xu, et al. Lasot: A high-quality large-scale single object tracking benchmark. *International Journal of Computer Vision (IJCV)*, 129(2):439–461, 2021. [2](#), [6](#), [13](#), [15](#), [18](#)
- [19] Heng Fan, Liting Lin, Fan Yang, Peng Chu, Ge Deng, Sijia Yu, Hexin Bai, Yong Xu, Chunyuan Liao, and Haibin Ling. Lasot: A high-quality benchmark for large-scale single object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [2](#), [3](#), [4](#), [6](#), [8](#), [14](#), [15](#)
- [20] J. Ferryman and A. Shahrokni. Pets2009: Dataset and challenge. In *2009 Twelfth IEEE International Workshop on Performance Evaluation of Tracking and Surveillance*, pages 1–6, 2009. [3](#)
- [21] Hamed Kiani Galoogahi, Ashton Fagg, Chen Huang, Deva Ramanan, and Simon Lucey. Need for speed: A benchmark for higher frame rate object tracking. In *ICCV*, 2017. [2](#), [6](#), [18](#)
- [22] Shenyan Gao, Chunlun Zhou, Chao Ma, Xinggang Wang, and Junsong Yuan. Aiatrack: Attention in attention for transformer visual tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 146–164, 2022. [8](#), [15](#)
- [23] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. [1](#), [2](#)
- [24] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [7](#)
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. [12](#)

- [26] João F. Henriques, Rui Caseiro, Pedro Martins, and Jorge Batista. High-speed tracking with kernelized correlation filters. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 37(3):583–596, 2015. 1
- [27] Lianghua Huang, Xin Zhao, and Kaiqi Huang. Globaltrack: A simple and strong baseline for long-term tracking. In *Proceedings of the Conference on Artificial Intelligence (AAAI)*, February 2020. 2
- [28] Lianghua Huang, Xin Zhao, and Kaiqi Huang. Got-10k: A large high-diversity benchmark for generic object tracking in the wild. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 43(5):1562–1577, 2021. 2, 6, 18
- [29] Mingzhen Huang, Xiaoxing Li, Jun Hu, Honghong Peng, and Siwei Lyu. Tracking multiple deformable objects in egocentric videos. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Vancouver, Canada, 2023. 2, 18
- [30] Matej Kristan, Aleš Leonardis, Jiří Matas, Michael Felsberg, Roman Pflugfelder, Joni-Kristian Kämäräinen, Martin Danelljan, Luka Čehovin Zajc, Alan Lukežič, Ondrej Drbohlav, Linbo He, Yushan Zhang, Song Yan, Jinyu Yang, Gustavo Fernández, and et al. The eighth visual object tracking vot2020 challenge results. In *Proceedings of the European Conference on Computer Vision Workshops (ECCVW)*, August 2020. 2
- [31] Matej Kristan, Ales Leonardis, Jiri Matas, Michael Felsberg, Roman Pflugfelder, Luka Cehovin Zajc, Tomas Vojir, Goutam Bhat, Alan Lukezic, Abdelrahman Eldesokey, Gustavo Fernandez, and et al. The sixth visual object tracking vot2018 challenge results. In *Proceedings of the European Conference on Computer Vision Workshops (ECCVW)*, September 2018. 6
- [32] Matej Kristan, Jiri Matas, Aleš Leonardis, Tomas Vojir, Roman Pflugfelder, Gustavo Fernandez, Georg Nebehay, Fatih Porikli, and Luka Čehovin. A novel performance evaluation methodology for single-target trackers. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 38(11):2137–2155, 2016. 1, 2, 3
- [33] Alina Kuznetsova, Aakrati Talati, Yiwen Luo, Keith Simmons, and Vittorio Ferrari. Efficient video annotation with visual interpolation and frame selection guidance. In *IEEE Winter Conference on Applications of Computer Vision, WACV 2021, Waikoloa, HI, USA, January 3-8, 2021*, 2021. 4
- [34] Bo Li, Wei Wu, Qiang Wang, Fangyi Zhang, Junliang Xing, and Junjie Yan. Siamrpn++: Evolution of siamese visual tracking with very deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 1, 8
- [35] Bo Li, Junjie Yan, Wei Wu, Zheng Zhu, and Xiaolin Hu. High performance visual tracking with siamese region proposal network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 8
- [36] Siyuan Li, Tobias Fischer, Lei Ke, Henghui Ding, Martin Danelljan, and Fisher Yu. Ovtrack: Open-vocabulary multiple object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5567–5577, 2023. 2, 6, 7
- [37] Liting Lin, Heng Fan, Zhipeng Zhang, Yong Xu, and Haibin Ling. Swintrack: A simple and strong baseline for transformer tracking. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. 3, 8
- [38] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014. 6
- [39] Yang Liu, Idil Esen Zulfikar, Jonathon Luiten, Achal Dave, Deva Ramanan, Bastian Leibe, Aljoša Ošep, and Laura Leal-Taixé. Opening up open world tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19045–19055, June 2022. 2, 8, 13
- [40] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10012–10022, October 2021. 12
- [41] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019. 12
- [42] Jonathon Luiten, Aljosa Osep, Patrick Dendorfer, Philip H. S. Torr, Andreas Geiger, Laura Leal-Taixé, and Bastian Leibe. HOTA: A higher order metric for evaluating multi-object tracking. *International Journal of Computer Vision (IJCV)*, 129(2):548–578, 2021. 3
- [43] Alan Lukežič, Luka Čehovin Zajc, Tomáš Vojř, Jiří Matas, and Matej Kristan. Now you see me: evaluating performance in long-term visual tracking, 2018. 3, 6
- [44] Fan Ma, Mike Zheng Shou, Linchao Zhu, Haoqi Fan, Yilei Xu, Yi Yang, and Zhicheng Yan. Unified transformer tracker for object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8781–8790, June 2022. 3, 8
- [45] Christoph Mayer, Martin Danelljan, Goutam Bhat, Matthieu Paul, Danda Pani Paudel, Fisher Yu, and Luc Van Gool. Transforming model prediction for tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8731–8740, June 2022. 2, 3, 4, 6, 8, 12, 15
- [46] Christoph Mayer, Martin Danelljan, Danda Pani Paudel, and Luc Van Gool. Learning target candidate association to keep track of what not to track. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 13444–13454, October 2021. 6, 8, 15
- [47] Matthias Mueller, Neil Smith, and Bernard Ghanem. A benchmark and simulator for uav tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, October 2016. 2, 6, 15, 18
- [48] Matthias Müller, Adel Bibi, Silvio Giancola, Salman Al-Subaihi, and Bernard Ghanem. Trackingnet: A large-scale

- dataset and benchmark for object tracking in the wild. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. 2, 6, 8, 18
- [49] Mubashir Noman, Wafa H Al Ghallabi, Daniya Kareem, Christoph Mayer, Akshay Dudhane, Martin Danelljan, Hisham Cholakkal, Salman Khan, Luc Van Gool, and Fahad Shahbaz Khan. Avist: A benchmark for visual object tracking in adverse visibility. In *33rd British Machine Vision Conference BMVC*, 2022. 6, 13, 14, 15
- [50] Jiangmiao Pang, Linlu Qiu, Xia Li, Haofeng Chen, Qi Li, Trevor Darrell, and Fisher Yu. Quasi-dense similarity learning for multiple object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 164–173, June 2021. 6, 7
- [51] Matthieu Paul, Martin Danelljan, Christoph Mayer, and Luc Van Gool. Robust visual tracking by segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 571–588, 2022. 8
- [52] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [53] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbelaez, Alexander Sorkine-Hornung, and Luc Van Gool. The 2017 DAVIS challenge on video object segmentation. *CoRR*, abs/1704.00675, 2017. 18
- [54] Hamid Rezaatfighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 6, 12
- [55] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 6, 13, 15, 17
- [56] Siyu Tang, Mykhaylo Andriluka, Bjoern Andres, and Bernt Schiele. Multiple people tracking by lifted multicut and person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 1, 2
- [57] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. FCOS: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 5
- [58] Jack Valmadre, Luca Bertinetto, João F. Henriques, Ran Tao, Andrea Vedaldi, Arnold W.M. Smeulders, Philip H.S. Torr, and Efstratios Gavves. Long-term tracking in the wild: a benchmark. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018. 2, 4, 18
- [59] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 3
- [60] Paul Voigtlaender, Jonathon Luiten, Philip H.S. Torr, and Bastian Leibe. Siam R-CNN: Visual tracking by re-detection. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2
- [61] Yi Wu, Jongwoo Lim, and Ming-Hsuan Yang. Object tracking benchmark. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 37(9):1834–1848, 2015. 1, 2, 3, 6, 18
- [62] Ning Xu, Linjie Yang, Yuchen Fan, Dingcheng Yue, Yuchen Liang, Jianchao Yang, and Thomas Huang. Youtube-vos: A large-scale video object segmentation benchmark, 2018. 2, 6, 18
- [63] Yinda Xu, Zeyu Wang, Zuoxin Li, Ye Yuan, and Gang Yu. Siamfc++: Towards robust and accurate visual tracking with target estimation guidelines. In *Proceedings of the Conference on Artificial Intelligence (AAAI)*, February 2020. 5
- [64] Bin Yan, Yi Jiang, Peize Sun, Dong Wang, Zehuan Yuan, Ping Luo, and Huchuan Lu. Towards grand unification of object tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 733–751. Springer International Publishing, 2022. 3, 8
- [65] Bin Yan, Houwen Peng, Jianlong Fu, Dong Wang, and Huchuan Lu. Learning spatio-temporal transformer for visual tracking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10448–10457, October 2021. 2, 3, 6, 8
- [66] Bin Yan, Haojie Zhao, Dong Wang, Huchuan Lu, and Xiaoyun Yang. 'skimming-perusal' tracking: A framework for real-time and robust long-term tracking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 2
- [67] Botao Ye, Hong Chang, Bingpeng Ma, Shiguang Shan, and Xilin Chen. Joint feature learning and relation modeling for tracking: A one-stream framework. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 341–357. Springer Nature Switzerland, 2022. 3, 8, 15
- [68] Bin Yu, Ming Tang, Linyu Zheng, Guibo Zhu, Jinqiao Wang, Hao Feng, Xuetao Feng, and Hanqing Lu. High-performance discriminative tracking with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9856–9865, October 2021. 3
- [69] 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 (CVPR)*, June 2020. 1, 2, 18
- [70] Qian Yu, Gérard Medioni, and Isaac Cohen. Multiple target tracking using spatio-temporal markov chain monte carlo data association. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2007. 1
- [71] Li Zhang, Yuan Li, and Ramakant Nevatia. Global data association for multi-object tracking using network flows. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2008. 1
- [72] Yifu Zhang, Peize Sun, Yi Jiang, Dongdong Yu, Fucheng Weng, Zehuan Yuan, Ping Luo, Wenyu Liu, and Xinggang

Wang. Bytetrack: Multi-object tracking by associating every detection box. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 1–21, 2022. 13, 14

- [73] Zikun Zhou, Jianqiu Chen, Wenjie Pei, Kaige Mao, Hongpeng Wang, and Zhenyu He. Global tracking via ensemble of local trackers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8761–8770, June 2022. 15

Appendices

In this supplementary material, we first give an overview of the different task definitions and the corresponding abbreviations used in the paper and supplementary material. Next, we describe the details of the model architecture and training in Sec. B. Then, we provide more insights into the experiments presented in the main paper and provide additional results on less popular tracking datasets in Sec. C. Then, we show visual results between the baseline and our tracker on multiple sequences of the proposed datasets including failure cases D. Next, we discuss the limitations of the proposed tracker in Sec. E. Finally, we provide additional insights about our dataset and compare it to datasets of related tasks in Sec. F.

A. Glossary

In this Section we will briefly summarize the different task definitions behind the individual abbreviations:

GOT. Generic Object Tracking refers to the task of tracking potentially multiple user-defined target objects of arbitrary classes specified by a user-specified bounding box in the initial video frame.

SOT. Single Object Tracking is the same task as GOT but focuses on the setting where only a single generic object needs to be tracked.

Multi-Object GOT. The same as GOT but emphasizes that multiple-objects need to be tracked. We use multi-object GOT because GOT is in other research works sometimes used interchangeably with SOT.

MOT. Multi Object Tracking is completely different from the tasks listed above because it requires a class category list to detect and track all objects corresponding to the defined class categories.

GMOT. Generic Multi Object Tracking is the same as MOT but instead of using a class category list to define the target objects, a single user-specified box shows an example object of the target class category. Thus, all objects that belong to the same class as the user-specified example need to be detected and tracked.

B. Model Architecture and Training Details

Architecture. We extract backbone features either from the ResNet-50 or the SwinBase backbone. For both backbones we extract the features corresponding to the blocks with stride 8 and 16. We only use the features with stride 16 for object encoding and feed these features into the model predictor. For both backbones we use a linear layer to decrease the number of channels from 1024 to 256 or 512 to 256 respectively. Thus we use 256 dimensional object embeddings e_i and a MLP to project the LTRB bounding box encoding map from 4 to 256 channels. Since the model predictor produces 256 dimensional convolutional filters we require the same number of channels for the FPN output features. In particular we use a two layer FPN that uses as input the enhanced Transformer encoder output features corresponding to the test frame as well as the aforementioned high resolution backbone test features. The high resolution input features have either 512 or 256 channels for the Resnet-50 or the SwinBase backbone respectively. Thus, we adapt the FPN accordingly depending on the used backbone.

Training Details. Since our tracker operates on full frames, we retain the full training and testing frames. The frames are re-scaled and padded to a resolution of 384×576 . As we use the feature maps with stride 16 for both the ResNet-50 [25] and SwinBase [40] backbones, this results in an extracted feature and score map resolution of 24×36 . For ResNet-50 we use pretrained weights on ImageNet-1k and for SwinBase on ImageNet-22k. We use a fixed size Gaussian when producing the score map encoding for each object where $\sigma = 0.25$. Furthermore, we use gradient norm clipping with the parameter 0.1 in order to stabilize training. In addition, we employ data augmentation techniques during training such as random scaling and cropping in addition to color jittering and randomly flipping the frame. The regression loss is given by

$$L_{\text{bbreg}} = \sum_{i=0}^n L_{\text{GIoU}} \left(\hat{b}_i^{\text{ltrb}}, \hat{b}_i^{\text{trb}} \right), \quad (5)$$

where L_{GIoU} denotes the generalized IoU-Loss [54]. The overall training loss is then defined as

$$L_{\text{tot}} = \lambda_{\text{cls}} L_{\text{cls}}(\hat{y}, y) + \lambda_{\text{bbreg}} \cdot L_{\text{bbreg}}(\hat{b}^{\text{ltrb}}, b^{\text{ltrb}}) \quad (6)$$

where $\lambda_{\text{cls}} = 100$ and $\lambda_{\text{bbreg}} = 1$ are scalars weighting the contribution of each loss component. We use ADAMW [41] with a learning rate of 0.0001 that we decay after 150 and 250 epochs by a factor of 0.2 and train all models on four A100 GPUs with a batch size of 4×12 or 4×6 .

Inference. During inference we adopt the simple memory updating approach described in [45]. In particular, updating the memory refers to adding a second dynamic training frame using predicted box annotations. We replace the

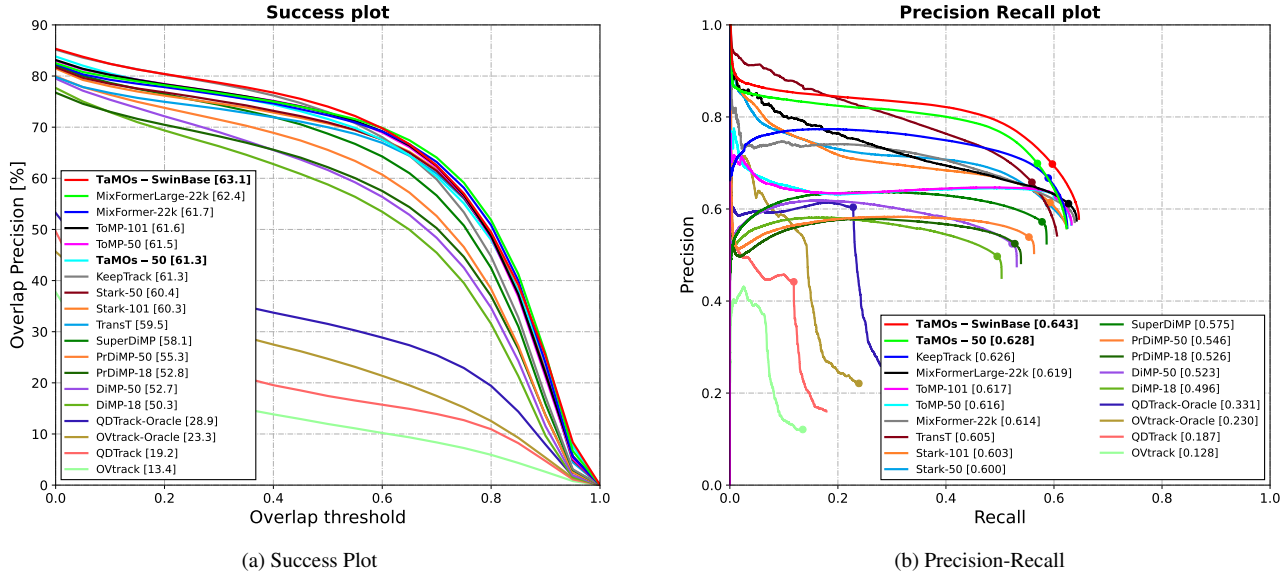


Figure 5. Success plot, showing OP_T , on LaGOT (AUC is reported in the legend). Tracking Precision-Recall curve on LaGOT – VOTLT is reported in the legend (the highest F1-score).

Table 7. Comparison of the combination of GOT and MOT methods. GOT return the detections and the MOT methods are used for object association over time on LaGOT.

GOT	MOT	F1-Score	Success	HOTA	MOTA	IDF1
TaMOS-SwinBase	—	0.643	63.1	62.1	58.2	74.7
	SORT [3]	0.438	35.7	45.9	52.2	43.3
	ByteTrack [72]	0.459	37.7	50.4	57.1	53.9
MixFormerLarge-22k	—	0.619	62.4	61.5	52.3	74.3
	SORT [3]	0.418	34.0	45.6	43.9	44.9
	ByteTrack [72]	0.450	36.4	47.5	44.8	49.6

second training frame (update the memory) if the maximal value in each score map is above the threshold of $\tau = 0.85$.

For accurate bounding box prediction and localization we employed an FPN. In contrast to training, where we applied the target models directly on the Transformer encoder features and also on the low- and high-resolution FPN feature maps, we only use the high-resolution score and bounding box prediction maps during inference. We empirically observed better training performance when applying the losses on each instead of only on the high resolution outputs. However, during inference we are only interested in the high resolution predictions.

C. Experiments

We provide more detailed results to complement the comparison shown in the main paper. In addition we provide result for the LaSOText [18] dataset in order to assess the performance of our tracker on sequences containing small objects. Similarly, we analyze the capability of our

tracker to handle adverse tracking conditions on AViT [49]. Furthermore, to provide results on another multiple object dataset we run the tracker on ImageNetVID [55].

C.1. LaGOT

To complement the results shown in the main paper, we report in Fig. 5 and Tab. 8 results for additional trackers and different variants, such as using a different backbone or different hyper-parameters. In Tab. 8 we report additional MOT sub-metrics and statistics on LaGOT. In general we conclude, that using larger backbones especially if they are pretrained on ImageNet-22k leads to the best results. Furthermore, we observe that the MOT methods QDTrack and OVTrack (evaluated with default parameters provided in the OVTrack GitHub repository¹) are not competitive with GOT methods. In particular, we observe that QDTrack and OVTrack achieve very low OWTA scores that depend on the Detection Recall (DetRe) and the Association Accuracy (AssA) scores. OVTrack scores the lowest DetRe despite being an open-vocabulary detector. While this is an expected limitation, we further observe that QDTrack achieves by far the lowest AssA caused by the poor Association Recall (AssRe) of 30.3 compared to DiMP-18 that achieves 67.6.

In addition to the SOT and MOT baselines presented in the main paper, we also evaluate an open-world tracker [39]. Such a tracker aims at tracking all objects in the scene and should therefore also be able to track the generic objects contained in LaGOT. In particular, we follow Liu *et al.* [39]

¹<https://github.com/SysCV/ovtrack>

Table 8. Comparison of different trackers using MOT metrics on LaGOT.

		HOTA	DetA	AssA	DetRe	DetPr	AssRe	AssPr	LocA	OWTA	MOTA	IDSW	IDF1
GOT	TaMOs-SwinBase	62.1	57.3	68.4	69.9	69.9	75.9	75.9	84.2	68.9	58.2	6734	74.7
	TaMOs-50	60.0	54.6	66.9	67.7	67.7	74.5	74.5	84.0	67.1	52.9	7901	72.0
SOT	MixformerLarge-22k	61.5	53.8	70.9	67.4	67.4	77.8	77.8	84.8	69.0	52.3	3150	74.3
	Mixformer-22k	61.2	54.0	70.0	67.4	67.4	77.0	77.0	84.5	68.6	53.2	3339	74.4
	ToMP-101	60.1	53.0	68.8	66.4	66.4	76.2	76.2	83.9	67.5	51.9	2638	73.8
	ToMP-50	60.0	53.0	68.6	66.4	66.4	76.0	76.1	83.8	67.4	52.3	2378	74.0
	STARK-ST-101	59.4	51.8	68.8	65.6	65.6	75.9	75.9	84.2	67.1	49.0	3568	72.5
	STARK-ST-50	59.4	51.9	68.5	65.6	65.6	75.6	75.6	83.9	66.9	49.5	4277	72.6
	TransT	57.8	50.2	67.1	64.3	64.3	74.5	74.6	84.3	65.6	46.6	2323	70.7
	KeepTrack	59.1	52.3	67.3	65.4	65.4	74.7	74.7	82.3	66.2	51.3	2299	73.8
	SuperDiMP	56.1	48.3	65.8	62.1	62.1	73.5	73.5	82.2	63.8	43.2	1966	69.7
	PrDiMP-50	53.0	45.6	62.1	59.6	59.6	70.3	70.4	81.3	60.7	38.4	2380	66.6
	PrDiMP-18	51.4	42.8	62.2	57.2	57.2	70.2	70.3	81.3	59.5	31.9	1981	63.4
	DiMP-50	50.8	42.1	62.0	56.2	56.2	69.7	69.7	80.2	58.9	29.4	1680	62.1
	DiMP-18	48.1	39.3	59.6	53.5	53.5	67.6	67.6	79.5	56.3	23.2	1757	59.0
	MOT	QDTrack	22.2	17.3	29.0	46.2	21.0	30.3	80.0	81.8	36.3	-115.8	18521
OVTrack		24.4	20.3	29.9	22.7	59.7	31.2	78.2	82.0	25.9	13.9	4951	23.5

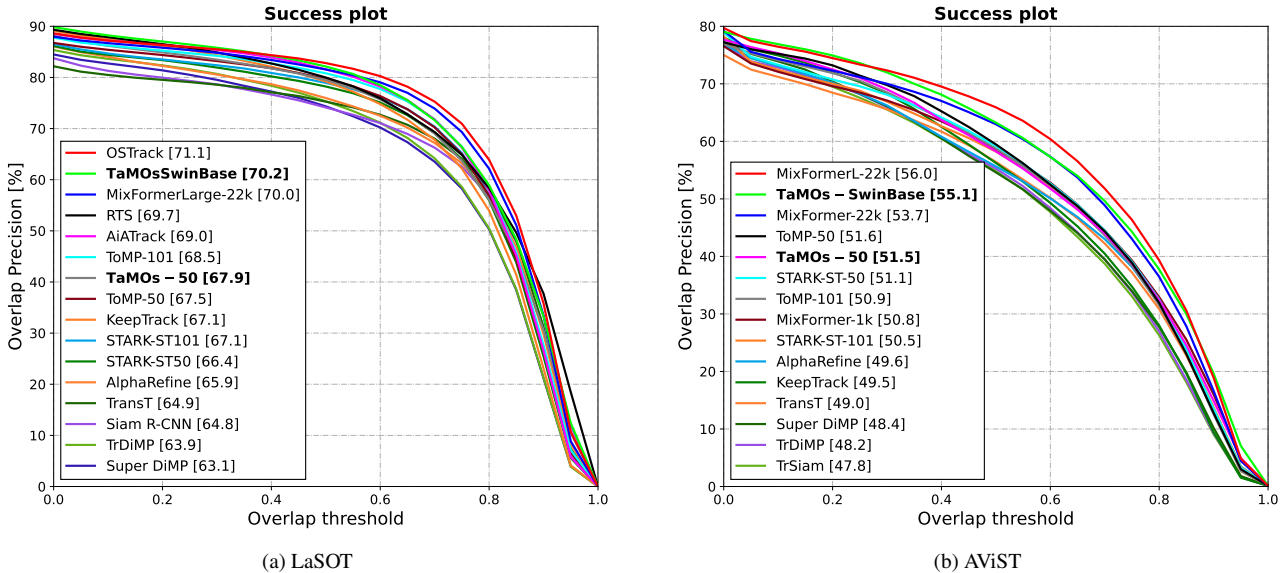


Figure 6. Success plot, showing OP_T , on LaSOT [19] and AViST [49] (AUC is reported in the legend).

and generate object proposals for each video frame using their provided open-world detector. Then, we run SORT [3] on top of the generated proposals using the default parameters. This leads to an OWTA score of 12.58, AssA of 3.57 and DetRe of 46.72. We conclude that the complex videos with long tracks of the proposed benchmark are for now too challenging for existing open-world trackers.

Finally, we add another experiment where we use our tracker TaMOs or the SOT tracker Mixformer as one-shot object detectors and feed their detections and scores to a MOT tracker that focuses on building the final tracklets. In particular we use the popular SORT [3] tracker and the recent state-of-the-art tracker ByteTrack [72]. For TaMOs and

Mixformer, using their predicted bounding boxes and object ids leads to far better results than using an MOT method on top for post-processing. This behaviour holds when measuring the performance of the resulting trackers with GOT as well as with MOT metrics, see Tab. 7. While there is potential to increase the robustness of GOT trackers in case of multiple objects, directly applying MOT trackers is not a good solution. Instead dedicated association algorithms for multi-object GOT are needed. We conclude, that TaMOs and the proposed SOT trackers run in parallel, are solid baselines for LaGOT.

Table 9. Comparison to the state of the art on LaSOTExt [18].

Method	Venue	LaSOTExt [18]		
		Prec	N-Prec	Succ
TaMOs-SwinBase		58.0	57.8	49.2
TaMOs-Resnet-50		54.1	55.0	46.7
AiATrack [22]	ECCV'22	54.7	58.8	49.0
OSTrack [67]	ECCV'22	57.6	61.3	50.5
ToMP-101 [45]	CVPR'22	52.6	58.1	45.9
ToMP-50 [45]	CVPR'22	51.9	57.6	45.4
GTELT [73]	CVPR'22	52.4	54.2	45.0
KeepTrack [46]	ICCV'21	54.7	61.7	48.2
SuperDiMP [14]	CVPR'20	49.0	56.3	43.7
LTMU [10]	CVPR'20	45.4	53.6	41.4
DiMP [4]	ICCV'19	43.2	49.6	39.2
ATOM [11]	CVPR'19	41.2	49.6	37.6

Table 10. Analysis of the FPN and the zooming mechanism on LaSOTExt [18] and UAV123 [47].

Backbone	FPN	Zoom	LaSOTExt	UAV123
			AUC	AUC
Resnet-50	×	×	41.3	56.2
Resnet-50	✓	×	43.1	58.2
Resnet-50	✓	✓	46.7	64.2
SwinBase	×	×	43.9	56.5
SwinBase	✓	×	44.6	57.3
SwinBase	✓	✓	49.2	66.2

C.2. LaSOT

In addition to the result table, shown in the main paper, we show in Fig. 6a the success plot for LaSOT [19]. We observe that our tracker is the most robust ($T < 0.3$). Furthermore, the plot shows that both MixFormerLarge-22k and OSTRack can regress more accurate bounding boxes ($0.5 < T < 0.9$). However, unlike these specialized single-target object trackers, our approach is capable of jointly tracking multiple targets.

C.3. LaSOTExt

Since our tracker always operates on the full frame without the help of a local search region, tracking small objects is challenging. Thus, we integrated an FPN in our tracker to improve the tracking accuracy. To analyze our tracker on small objects we run it on LaSOTExt [18] and UAV123 [47]. Tab. 10 shows that including an FPN improves the tracking results on both datasets but is more effective when using a Resnet-50 as backbone.

To track small objects a high feature map resolution is desirable. To better cope with extremely small objects, found in some SOT benchmarks, we add a simple zooming mechanism. In particular, when the target is smaller than 30×30 pixels, we crop a region of the image that en-

ures this minimal target size when up-scaled to the input-resolution of 384×576 . Tab. 10 clearly shows that using such a zooming mechanism improves the results on LaSOTExt and UAV123 considerably, due to the presence of extremely small objects in these datasets.

Tab. 9 shows that our tracker with FPN and zooming achieves competitive results on LaSOTExt. In particular it achieves the highest precision and the second highest success AUC only being outperformed by OSTRack [67].

C.4. AViT

In order to validate our tracker in adverse visibility scenarios we run it on AViT [49]. Fig 6b shows that our tracker achieves excellent results with a success AUC of 55.1. This result shows that our tracker is able to track generic single objects even in visually challenging scenarios. The best tracker MixFormerLarge-22k is able to regress more accurate bounding boxes ($0.3 < T < 0.9$), as it relies on small search area selection to ensure high-resolution features. In contrast, our approach is capable of jointly tracking multiple objects.

C.5. ImageNetVID

In order to validate the proposed multiple object GOT tracker not only on LaGOT but also on another multiple object dataset, we modify ImageNetVID [55]. Since ImageNetVID is a video object detection datasets instead of a GOT dataset we perform the following adaptations. First, we remove all tracks that are not present in the first frame. Then, we use the remaining tracks to produce the bounding box annotations of the first frame. For simplicity we remove the 11 sequence where no track is visible in the first frame. This results in 544 sequences with 938 tracks and 1.7 tracks on average per video. Fig. 8 shows the success plot on the resulting multiple object GOT dataset. We observe that all trackers achieve relatively high AUC mostly differing in bounding box accuracy. Both versions of our tracker outperform the baselines ToMP-50 and ToMP-101 [45]. In particular, we notice the superior bounding box accuracy of our tracker compared to ToMP. To summarize we observe a similar ranking between trackers on ImageNetVID and the proposed LaGOT dataset. However, LaGOT is more challenging due to the higher average track number (2.9 vs. 1.7) and the much longer sequence length (2258 vs. 312) that leads more frequently to occlusions and out-of-view events.

D. Visual Results

Visual Comparison to the State of the Art. We show visualizations of the tracking results of the baseline (ToMP-101) and our proposed tracker (TaMOs-SwinBase) on four different sequences of the proposed LaGOT benchmark in Fig. 7. The first frame specifies the target objects annotated

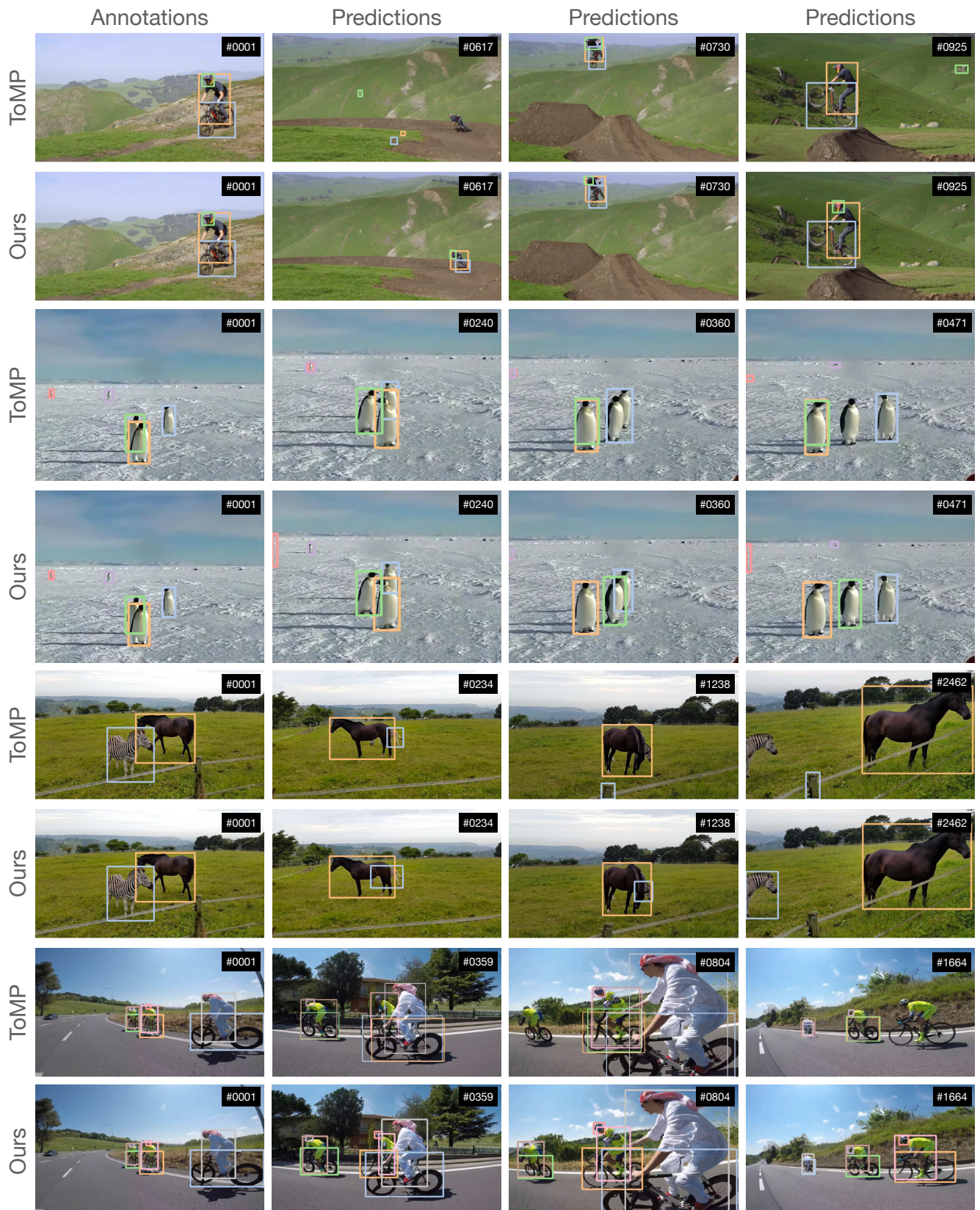


Figure 7. Visual comparison between the proposed tracker (Ours-SwinBase) and the baseline ToMP-101 on different LaGOT sequences.

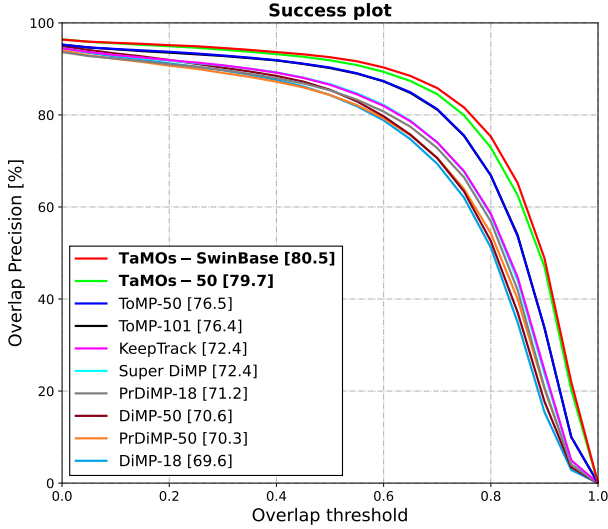


Figure 8. Success plot, showing OP_T , on ImagenetVID [55] (AUC is reported in the legend).

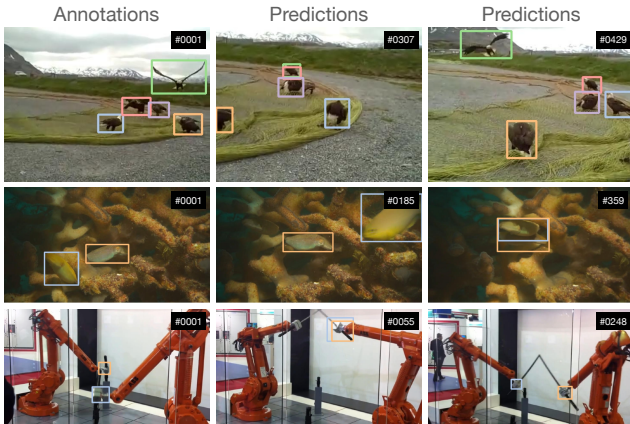


Figure 9. Visual examples of failure cases of the proposed tracker (Ours-SwinBase) on different LaGOT sequences.

with bounding boxes that should be tracked in the video. The other frames show predictions of both trackers. The results on the first and third sequences demonstrate that our tracker can re-detect occluded objects quickly whereas a search area based tracker is not able to re-detect the targets if they reappear outside of the search area. The second and fourth sequences show the superior robustness of our tracker. It is able to distinguish similarly looking objects better without confusing their ids. For more visual results we refer the reader to the mp4-videos submitted alongside this document. Each video shows the predictions of the proposed tracker TaMOs-SwinBase on the proposed LaGOT benchmark. Please note that we always produce a bounding box for visualization independent of its confidence score.

Failure Cases. Fig. 9 shows typical failure cases of the pro-

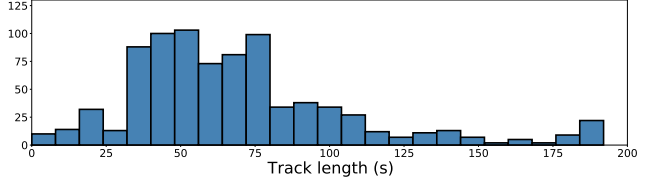


Figure 10. Track lengths distribution of the LaGOT benchmark.

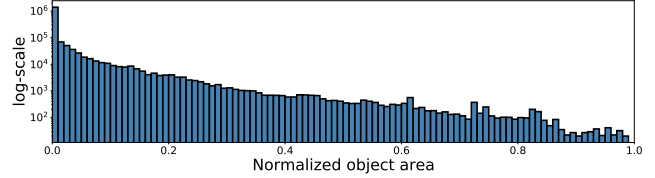


Figure 11. Object size distribution of the LaGOT benchmark.

posed tracker on three different sequences of the proposed LaGOT benchmark. Particularly challenging are videos that contain multiple visually similar objects since our tracker does not employ any motion model but rather tracks the objects via the learned appearance from the first frame. Another failure case occurs when the target object is no longer visible such that our tracker might start to track a visually similar distractor instead. However, once the target reappears our globally operating tracker is usually able to re-detect it. Lastly, if multiple visually similar objects need to be tracked our tracker might fail to distinguish these objects such that it produces multiple bounding boxes with different ids for the same object.

E. Limitations and Future Work

Currently the number of objects that can be tracked is limited by the pool-size of the object embeddings. While it is possible to learn a larger pool-size it is cumbersome. Thus, an interesting direction for future research would be to generate an arbitrary number of object embedding on the fly such that any number of target objects can be tracked.

Furthermore, we propose to use an FPN to regress more accurate bounding boxes for small objects and show that adding such an FPN helps. However, as in object detection, tracking extremely small objects is challenging due to the limited feature resolution when processing the full frame.

F. Datasets

Below we provide additional details about our annotated dataset, such as examples of new classes and various statistics, as well as an extensive comparison to existing datasets that focus on related tasks.

Table 11. Comparison of LaGOT and the existing datasets. Statistics is provided for test or validation set for the datasets for which test set annotations are hidden. * For MOT15-20 we report stats on the train set.

Dataset	Num Classes	Num Videos	Avg Video length (num frames)	Avg Tracks per Video	Avg Track Length (num boxes)	Avg Track Length (s)	Avg Instances per frame	Video FPS	Annotation FPS
YouTubeVOS [62]	91	474	135	1.74	27	4.5	1.64	30 FPS	6 FPS
Davis17 [53]	-	30	67	1.97	67	2.8	1.97	24 FPS	24 FPS
ImageNetVID* [17]	30	555	317	2.35	208	7	1.58	30 FPS	30 FPS
TAO* [15]	302	988	1010	5.55	21	21	3.31	30 FPS	1 FPS
BDD100k [69]	11	200	198	94.21	26	5	11.8	30 FPS	5 FPS
MOT15 [16]*	1	11	500	45.5	75	3	8	2.5-30 FPS	2.5-30 FPS
MOT16 [16]*	1	7	760	74	273	10	38	14-30 FPS	14-30 FPS
MOT20 [16]*	1	4	2233	583	572	23	150	25 FPS	25 FPS
DogThruGlasses [29]	1	30	419	3.3	352.6	11.7	2.4	30 FPS	30 FPS
GMOT-40 [1]	10	40	240	50.65	133	5.3	26.6	24-30 FPS	24-30 FPS
TrackingNet [48]	27	511	442	1	442	15	1	30 FPS	30 FPS
UAV123 [47]	8	123	915	1	915	28	1	30 FPS	30 FPS
OTB-100 [61]	16	100	590	1	590	20	1	30 FPS	30 FPS
NFS-30 [21]	15	100	479	1	479	14	1	30 FPS	30 FPS
GOT10k [28]	84	420	150	1	150	15	1	10 FPS	10 FPS
OxUvA [58]	8	200	4198	1	60	140	1	30 FPS	1 FPS
LaSOT [18]	71	280	2430	1	2430	81	1	30 FPS	30 FPS
LaGOT	102	294	2258	2.89	707	71	2.41	30 FPS	10 FPS

F.1. Insights

Fig. 10 shows the distribution of the track lengths in seconds for all tracks in the proposed benchmark LaGOT. We observe that most tracks are between 30 and 110 seconds long. Furthermore, Fig. 11 shows the size distribution of the annotated objects in the dataset. We conclude that various sizes are present in the dataset but large objects are rare than small ones. Further, the distribution shows that the targets are not visible in a large amount of video frames indicated by an object area of zero.

During the annotation process, we added 31 new classes: *rotor, fish, backpack, motor, wheel, garbage, drum, accordion, super-mario, hockey puck, hockey stick, kite-tail, ball, crown, stick, spiderweb, head, banner, face, bench, tissue-bag, para glider, star-patch, shadow, bucket, helicopter, sonic, hero, ninja-turtle, reflection, rider.*

F.2. Comparison

We provide a detailed comparison of related existing datasets in Tab. 11. We divide the table into Video Object Segmentation (VOS), Video Object Detection, Multiple Object Tracking (MOT), Generic Multiple Object Tracking (GMOT) and Single Object Tracking (SOT) datasets.

The length of VOS sequences is much shorter than in our LaGOT benchmark (2.8s/4.5s vs 71s). Similarly the video object detection dataset ImagenetVID contains shorter sequences (7s vs. 71s), fewer classes (30 vs 102) and a smaller number of average tracks per sequence (2.35 vs 2.89) than LaGOT. MOT datasets typically focus on fewer classes, contain shorter sequences or are annotated at low frame

rates only. TAO contains many more classes than typical MOT datasets but provides annotations only at 1 FPS leading to a much lower average number of annotated frames per track than LaGOT (21 vs. 707). The GMOT-40 dataset contains fewer classes, fewer videos, shorter sequences and provides due to its task only annotations of one particular object class per sequence compared to LaGOT. In contrast to SOT datasets that provide only a single annotated object per sequence, LaGOT provides on average 2.89 tracks per sequence. Furthermore, it contains longer sequences than most listed SOT datasets. Overall LaGOT enables to properly evaluate the robustness and accuracy of multiple object GOT methods. A key factor are the multiple annotated tracks per sequence at a high frame rate and the relatively long sequences.