# Motion Prediction in Visual Object Tracking

Jianren Wang*, Yihui He*

*Abstract*— **Visual object tracking (VOT) is an essential component for many applications, such as autonomous driving or assistive robotics. However, recent works tend to develop accurate systems based on more computationally expensive feature extractors for better instance matching. In contrast, this work addresses the importance of motion prediction in VOT. We use an off-the-shelf object detector to obtain instance bounding boxes. Then, a combination of camera motion decouple and Kalman filter is used for state estimation. Although our baseline system is a straightforward combination of standard methods, we obtain state-of-the-art results. Our method establishes new state-of-the-art performance on VOT (VOT-2016 and VOT-2018). Our proposed method improves the EAO on VOT-2016 from 0.472 of prior art to 0.505, from 0.410 to 0.431 on VOT-2018. To show the generalizability, we also test our method on video object segmentation (VOS: DAVIS-2016 and DAVIS-2017) and observe consistent improvement.**

## I. INTRODUCTION

Tracking moving objects over space and time is fundamental for understanding the dynamic visual world, which has many practical applications in video processing, such as self-driving [2], video surveillance [3], and UAV navigation [4].

Many attempts have been addressed to improve the performance of trackers over the years. In the early days, motion model was a core component of tracking - constant velocity models [5], [6], Kalman filters [7], [8], particle filters [9], [10] and even social force models [11], [12] for more complicated motions. In fact, in the early days, it was the dominant component, because (i) decent appearance descriptors were not available, and (ii) it had its roots outside of computer vision (*e.g.,*, tracking point targets in RADAR data), where there is no appearance information. However, most modern trackers assuming zero-velocity model, because (i) Modern tracking datasets contain many sequences with random camera motion, which fails most motion model [13], [14] (ii) With better feature extraction and bounding box regression ability introduced by CNN [15], [16], modern trackers [17], [18] can rely less on motion priors, which is known as tracking by detection.

In contrast to prior works which tend to develop accurate systems based on more computational costly feature extractors, this work aims to develop a robust motion prediction method. We address the importance of motion prediction, even if trackers based only on appearance cues have already achieved good performance. We prove that even modern CNN trackers can benefit a lot from accurate motion predictions.

Concretely, in the first stage, we decouple camera motion and object motion. Second, we predict the object state in the future frame and create an adaptive search region for the detector to process. The adaptive search region focuses on smaller local regions when objects have slower speeds and smaller sizes, and vice versa. We then project the predicted state and search region back to the camera coordinate of the frame. We finally update the object state based on the measurement from the off-the-shelf object detector. Both state prediction and update are based on Kalman filter.

Our method has several benefits: First, by decoupling object motion from camera motion, we alleviate the motion noise caused by camera shake. Second, we free modern trackers from using only appearance information. As most tracking and segmentation methods can only discriminate foreground from the non-semantic background [19], the performance suffers significantly when the target object is surrounded by similar objects (know as distractors [19]). Our method can also improve the performance under occlusions since the motion model can prevent the detector from tracking the occluders. We show in the experiments that our method improves the tracking performance by a large margin under both cases. We visualize part of the results in Fig. 1. Third, we achieve robust bounding box prediction by updating the state through the Kalman filter.

We evaluate our framework on major tracking datasets: VOT-2016 [13] and VOT-2018 [14]. We demonstrate the effectiveness of our method, both qualitatively and quantitatively. On VOT-2016, we achieve 0.505 EAO, and on VOT-2018 we achieve 0.431 EAO. Although our method focuses on video object tracking, it generalizes well to video object segmentation. Consistent improvements over the baseline method are demonstrated on VOS (DAVIS-2016 [20] and DAVIS-2017 [21]).

We summarize our contributions as follows: First, we revisit motion prediction in visual object tracking, which has long been ignored. Second, we propose a method that combines motion decouple, motion prediction with off-the-shelf appearance-based trackers. Third, our proposed method achieves state-of-the-art performance on VOT and can also consistently improve the performance on VOS.

## II. RELATED WORKS

### A. Video Object Tracking

In tracking community, significant attention has been paid to discriminative correlation filters (DCF) based methods [22], [23], [24], [25]. These methods allow discriminating between the template of an arbitrary target and its 2D translations at a breakneck speed. MOSSE [22] is the
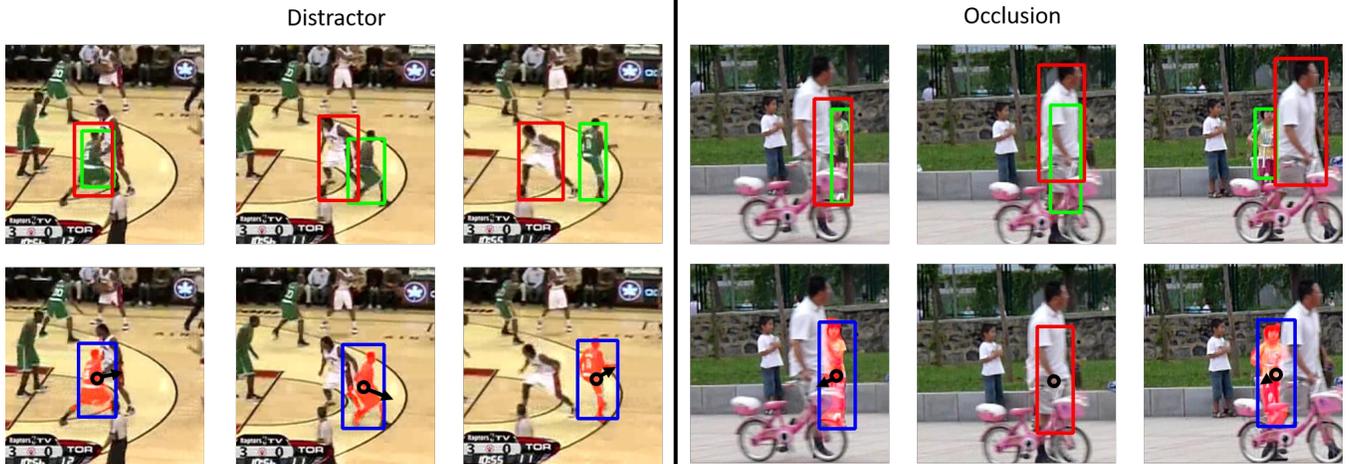
Fig. 1. The first row shows the results of state-of-the-art tracker SiamMask [1] (red) and ground truth (green). The second row shows the motion prediction (black arrow) and tracking results (with the blue bounding box and red segmentation mask) of our model. Our method improves the robustness against distractors and occlusions. (better view with color)

pioneering work which proposes a fast correlation tracker by minimizing the squared error. Performance of DCF-based trackers has then been notably improved through the using of multi-channel features [26], [27], [28], robust scale estimation [29], [30], reducing boundary effects [31], [32] and fusing multi-resolution features in the continuous spatial domain [33].

Tracking through Siamese Network is also an important approach [34], [35], [36], [37]. Instead of learning a discriminative classifier online, the idea is to train a deep siamese similarity function offline on pairs of video frames. At test time, this function is used to search for a candidate similar to the template given in the starting frame on a new video, once per frame. GOTURN [38] used a deep regression network to predict the motion between successive frames. SiamFC [36] implemented a fully convolutional network to produce a correlation response map with high values at target locations, which establishes a basic form of modern Siamese framework. Many following works have been proposed to improve the accuracy while maintaining fast inference speed by introducing semantic branch [39], region proposals [17], hard negative mining [19], ensembling [40], deeper backbone [18] or high-fidelity object representations [1].

Under the assumption that objects are under minor displacement and size change in consecutive frames, most modern trackers, including the ones mentioned above, use a steady search region, which is centered on the last estimated position of the target with the same ratio. Although it is very straightforward, this oversimplified prior often fails under occlusion, motion change, size change, or camera motion, as it is evident in the examples of Fig. 1. This motivates us to propose a robust motion prediction module that fits all these methods.

### B. Video Forecasting

The ability to predict and therefore to anticipate the future is an important attribute of intelligence. Many methods are proposed to improve the temporal stability of semantic video segmentation. Luc *et al.* [41] develop an auto-regressive convolutional neural network that learns to generate multiple future frames iteratively. Similarly, Walker *et al.* [42] use a VAE to model the possible future movements of humans in the pose space. Instead of generating future states directly, many methods attempt to propagate segmentation from preceding input frames [43], [44], [45].

Unlike previous work, we extract a motion model for each object and set up a new search region for detection and segmentation accordingly.

### III. METHOD

Our method first decouples object motion from camera motion. Second, we predict the object state in the future frame and create an adaptive search region for the detector to process. We then project the predicted state and search region back to the camera coordinate. We finally update the object state based on the measurement from the off-the-shelf object detector. State prediction and update are based on Kalman filter. We illustrate our framework in Fig. 2.

### A. Decouple

Object motion in a given image is the superposition of camera motion and object motion. These motions may lie in different modes (*e.g.,*, random camera shaking, or object moving direction sudden change). Thus, predicting object motion in camera coordinates for a long horizon will lead to instability. To solve this problem, we first pick a reference frame ($F_k$, $k$ denotes $k^{th}$ reference frame) every $n$ frames and thus separate the long video into several pieces of short n-frame videos.

Second, we adopt the method proposed by ARIT [46] to decouple the camera motion and object motion within each short video. ARIT assumes that pending detection frame ($F_{k+t}$) and its reference frame ($F_k$) are related by a homography ($H_{k,k+t}$). To estimate the homography, the
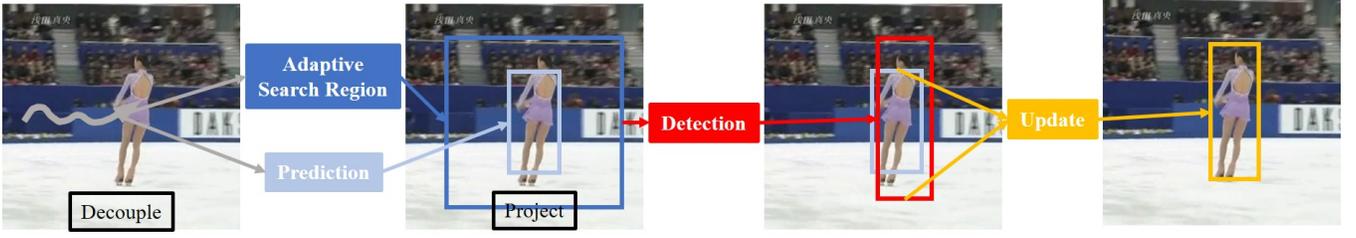
Fig. 2. Our method first decouple camera motion and object motion. Second, we predict the object state in the future frame and create an adaptive search region for the detector to process. We then project the predicted state and search region back to the camera coordinate. Finally, we update the object state based on the measurement from the off-the-shelf object detector. State prediction and update are based on Kalman filter.



Fig. 3. One example for decoupling camera motion and object motion (arrows illustrates the movement of object center).

first step is to find the correspondences between two frames. As mentioned in ARIT, we combine SURF features [47] and motion vectors from the optical flow to generate sufficient and complementary candidate matches, which is shown to be robust [48], [46]. Here we use PWCNet [49] for dense flow generation.

As a homography matrix contains eight free variables, at least four background points pairs should be used. We calculate the least square solution of eq. 1 and optimize it to obtain robust solution through RANSAC [50], where $p_k^{bp}$ and $p_{k+t}^{bp}$ denotes random selected background matching pairs in $F_k$ and $F_{k+t}$ using the above mentioned features.

$$H_{k,k+t} \times p_k^{bp} = p_{k+t}^{bp} \tag{1}$$

Fig. 3 illustrates the working principle of the decoupling step. The origin video for Fig. 3 is a handheld video with trembling background. The motion of the pedestrians in the origin video is highly unpredictable with huge background uncertainties. However, by mapping the target frame towards the reference frame, the movement for pedestrians could be more predictable and continuous.

There are two cases where motion decouple is not applied: (i) no correspondences between the future frame and the reference frame (ii) outlier of RANSAC is larger than an

error threshold. These always happen due to severe blur, where the tracking relies purely on target appearance.

For simplicity, the following calculations are under reference coordinate without further noticing.

B. Prediction

We parameterize the bounding box as a set of five parameters, including the coordinate of the object center $(x, y)$ and objects size $(w, h)$ and its confidence $c$. Our prediction does not need any training process and can be directly applied for inference. We formulate the state of object trajectory as a 6-dimensional vector $s = (x, y, w, h, v_x, v_y)$, where the additional variables $v_x, v_y$ represent the velocity of objects.

To predict the object state in the next frame, we use the dynamic model in the Kalman filter, which is shown in (eq. 2).

$$s_{t|t-1} = F s_{t-1|t-1} + B u_t \tag{2}$$

where $s_{t|t-1}$ is the prior state estimation given observations up to time $t-1$, $s_{t-1|t-1}$ is the optimal result of the previous state, $u_t$ is the control amount of the process at time t (zero in our case). And $F$ is the transition matrix, $B$ is the system parameter. In this paper, we approximate the inter-frame displacement of objects using the constant velocity model, which is initialized to zero for each object. Then we predict the covariance corresponding to the process result:

$$V_{t|t-1} = F V_{t-1|t-1} F^T + Q \tag{3}$$

where $V_{t|t-1}$ is a prediction of the covariance corresponding to the state $s_{t|t-1}$, $V_{t-1|t-1}$ is the covariance corresponding to the previous state $s_{t-1|t-1}$. And $Q$ is the covariance matrix of system noise (assumed to be Gaussian).

C. Adaptive Search Region

To alleviate the information needed to be processed by detectors and better filter out distractors, we dynamically set up a new search region in the coming frame centered at the predicted object position. The adaptive search region is modified with respect to the predicted velocity and object size.

Given the estimated position, we setup the search region (a $S \times S$ square) as following:

$$S = k\sqrt{(w + p)(h + p)} \tag{4}$$

$$k = 1 + 2 \times sigmoid(||v||_2 - \theta_v) \quad (5)$$

where $p = \dfrac{w + h}{2}$, $v$ is the predicted velocity and $\theta_v$ is a pre-defined threshold.

### D. Project

We then project the estimated adaptive search region and predicted state back to the future frame as following:

$$H_{k,k+t} \times P_k = P_{k+t} \quad (6)$$

where $P_k$ and $P_{k+t}$ are the key-points (centers and corners) in frame k and frame k+t. Since affine transformations do not respect lengths and angles, we recalculate object sizes based on projected corners.

### E. Detection

It is worth noticing that our method does not depend on specific detection or segmentation methods. In this paper, we adopt SiamRPN++ [18] for detection and SiamMask [1] for segmentation, since they achieve a good balance between accuracy and speed. We refer readers to [16], [17] for understanding the region proposal branch and [51], [52] for understanding the mask branch.

### F. Update

To account for the uncertainty in prediction, we update the entire state space of trajectory based on its corresponding measurement, i.e., the detection result, and obtain the final trajectories using the following equation:

$$s_{t|t} = s_{t|t-1} + K_t(D_t - Os_{t|t-1}) \quad (7)$$

where $D_t$ is the measurement (detection result in our case). And $O$ is the observation matrix, $K_t$ is the optimal Kalman gain defined by eq. 8. In our case, velocity states are not observable.

$$K_t = \dfrac{V_{t|t-1}O^T}{OV_{t|t-1}O^T + R} \quad (8)$$

where R is the covariance matrix corresponding to the measurement noise (assumed to be Gaussian).

We then perform the covariance update as following:

$$V_{t|t} = (I - K_tO)V_{t|t-1} \quad (9)$$

where $I$ is an identity matrix.

We refer reader to [7], [8] for more details on Kalman filter. We only execute the aforementioned update when the detection confidence score is larger than a threshold $\theta_d$. If the detection confidence score is less than $\theta_d$, we update object states use only eq. 2 and eq. 3. This can help to track objects under occlusions and large appearance changes.

The motion consistency between video frames in different sliced videos with different reference frames could be an issue because the initialization of the velocity for the reference frame could be critical to the accuracy of the position update. To maintain the motion consistency, we choose the $n_{th}$ frame, which is the last frame in the sliced video, as the next reference frame with the refined position and velocity estimation from Kalman filter based on the former reference frame. Therefore, the velocity of the object, with respect to the new frame, could be initialized by mapping the refined velocity towards the new reference.

## IV. EXPERIMENTS

In this section, we evaluate our approach on three tasks: motion prediction, visual object tracking (VOT-2016 and VOT-2018), and semi-supervised video object segmentation (DAVIS-2016 and DAVIS-2017).

### A. Evaluation of motion prediction

*a) Datasets and settings:* We use VOT-2016 [13] and VOT-2018 [14] to evaluate the performance of motion prediction. Both datasets contain 60 public sequences with different challenging factors: camera motion, object motion change, object size change, occlusion, and illumination change, which makes it extremely challenging for object motion prediction [14]. We use SiamMask [1] for detection, which returns a segmentation mask for each tracking object. We thus use the center of mass of predicted mask as detected object position $(x, y)$. Our prediction of position and velocity is calculated as mentioned in Section III-B. For the baseline, the predicted position of the next frame (t+1) is always the same as the current frame (t), while object velocity is always predicted as 0. The ground truth position is set as the center of the annotated rotated bounding box, while the velocity is the difference between two consecutive positions. We evaluate the position error from ground truth with Euclidean distance and velocity error with Euclidean distance, cosine distance, and magnitude distance. Cosine distance is the cosine value between predicted velocity and ground truth velocity (the higher, the better). Magnitude distance is the absolute difference between the absolute value of predicted velocity and ground truth velocity. We adopt the reinitialize mechanism as used in the official VOT toolkit. When the segmentation has no overlap with ground truth, we reinitialize the tracking method with ground truth after five frames.

*b) Results on VOT-2016 and VOT-2018:* Table.I presents the comparison of position prediction error using the baseline method and our method. As it is shown in the table, for both of these two datasets, our method could dramatically reduce the prediction errors of the object position. The mean square error for object position on VOT-2018 could be reduced by half from 16 pixels to 8 pixels. Meanwhile, Fig.4 shows when the object velocity is high, our method could provide a more accurate prediction compared with Baseline, which does not consider the influence of object motion. The results prove that the decoupling strategy could reduce the background uncertainty, and the Kalman filter would provide a relatively reliable prediction for object position in the next frame. Higher accuracy for object position prediction could benefit the generation of search regions for object

| Dataset | Tracker | Pos Err. |
|---------|---------|----------|
| VOT-2016 | Baseline | 16.281 |
| | Ours | 8.198 |
| VOT-2018 | Baseline | 14.593 |
| | Ours | 8.744 |

TABLE I

POSITION PREDICTION ERROR ON VOT-2016 AND VOT-2018



Fig. 4. Position predictions (red for Ours, yellow for Baseline and blue cross for ground truth) (better view with color)

tracking and eventually improve the performance of object segmentation.

For velocity, as can be seen in Table II, our method significantly reduce the estimation error. In VOT-2018, Ours achieves 0.763 cosine distance, which is about 37-degree divergence from ground truth velocity direction. The main cause of the error is that objects are not always rigid, thus "center of mass" can approximate the overall motion of the object (Fig. 5). The size change of objects will further increase the prediction error. However, with the correction procedure of the Kalman filter, this error (noise) can be stabilized. One possible solution to decrease velocity prediction error is tracking each part of non-rigid objects and grouping all parts together to get the final prediction [53].
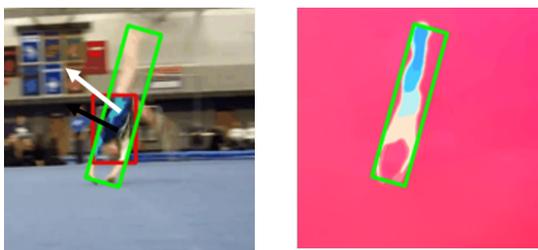


Fig. 5. Velocity predictions (Left) (white for goundtruth, black for prediction, both extended by 5 times longer for better visualization) Optical Flow (Right) (better view with color)

| Dataset | Tracker | MSE Err. | Cosine | Mag. |
|---------|---------|----------|--------|------|
| VOT-2016 | Baseline | 8.274 | - | 8.274 |
| | Ours | 4.596 | 0.667 | 3.190 |
| VOT-2018 | Baseline | 7.006 | - | 7.006 |
| | Ours | 4.298 | 0.793 | 2.929 |

TABLE II

VELOCITY PREDICTION ERROR ON VOT-2016 AND VOT-2018

| | VOT-2016 | | |
|---------|------|------|------|
| Trackers | A | R | EAO |
| SiamRPN++ | 0.633 | 0.181 | 0.472 |
| Ours | 0.642 | 0.139 | 0.505 |

TABLE III

COMPARISON WITH SIAMRPN++ ON VOT-2016

## B. Evaluation of VOT

We adopt two widely used benchmarks for the evaluation of the object tracking task: VOT-2016 and VOT-2018. Here we adopted SiamRPN++ as our detection module. We compare our method against the state-of-the-arts using the official metric: Expected Average Overlap (EAO), which considers both accuracy and robustness of a tracker [14]. We further conduct an experiment on VOT-2018 for evaluating the performance under different conditions.

*a) Results on VOT-2016:* Table III presents comparisons of tracking performance between our method and SiamRPN++ on VOT-2016 dataset. Our method improves the robustness by 23.2%, and provide a 7.0% gain of EAO, which achieves 0.505. The baseline is the state-of-the-art, other methods are not compared for simplicity.

*b) Results on VOT-2018:* In Table IV we compare our method against eleven recently published state-of-the-art trackers on the VOT-2018 benchmark (A stands for accuracy and R stands for robustness). We establish a new state-of-the-art tracker with 0.431 EAO and 0.607 accuracy. In particular, our model outperforms all existing Correlation Filter-based trackers. This is very easy to understand since our baseline SiamRPN++ relies on deeper feature extraction, which is much richer than all existing Correlation Filter-based methods. Interestingly, our method even outperforms the baseline method. Previous research shows Siamese based trackers have strong center bias despite the appearances of test targets [18]. Thus, by estimating the center of the search region more accurately, Siamese trackers can also achieve better regression result (*e.g.,*, bounding box detection, or object segmentation). Besides, our method achieves the lowest robustness among all Siamese based trackers. This is even exhilarating because one of the key vulnerability of Siamese based trackers is the low robustness. The main reason is that most Siamese networks can only discriminate foreground from the non-semantic background [19] and thus suffer from distinguishing target object from surrounding objects. Our proposed motion prediction module adopts a straightforward strategy and shows great improvement of robustness from 0.241 to 0.203, which provides another strategy to achieve better robustness: by setting more accurate and targeted search region. And our proposed modules only decrease the running speed by a small margin (5FPS) since all calculations are done in GPU (RTX2080Ti).

To further analysis where the improvements come from, we show the qualitative results of our method and the baseline SiamRPN++ (Fig. 6). Just as mentioned above, the robustness comes from less tracking object switching and

| | DaSiamRPN | SA_Siam_R | CPT | DeepSTSRCF | DRT | RCO | UPDT | SiamMask | SiamRPN | MFT | SiamRPN++ | Ours |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| EAO ↑ | 0.326 | 0.337 | 0.339 | 0.345 | 0.356 | 0.376 | 0.378 | 0.380 | 0.383 | 0.385 | 0.410 | **0.431** |
| Accuracy ↑ | 0.569 | 0.566 | 0.506 | 0.523 | 0.519 | 0.507 | 0.536 | 0.609 | 0.586 | 0.505 | 0.594 | 0.607 |
| Robustness ↓ | 0.337 | 0.258 | 0.239 | 0.215 | 0.201 | 0.155 | 0.184 | 0.276 | 0.276 | 0.140 | 0.241 | 0.203 |
| Speed (FPS) ↑ | 160 | 32 | 14 | 24 | < 1 | 7 | < 1 | 56 | 200 | < 1 | 35 | 30 |

TABLE IV

COMPARISON WITH THE STATE-OF-THE-ART TRACKERS UNDER EAO, ACCURACY, ROBUSTNESS AND SPEED ON THE VOT-2018 DATASET.

| Datasets | Methods | J | F |
|---|---|---|---|
| Davis-2016 | SiamMask | 0.713 | 0.674 |
| | Ours | 0.732 | 0.692 |
| Davis-2017 | SiamMask | 0.543 | 0.585 |
| | Ours | 0.554 | 0.604 |

TABLE V

J AND F RESULTS ON DAVIS-2016 AND DAVIS-2017

| | EAO | A | R |
|---|---|---|---|
| SiamMask | 0.380 | 0.609 | 0.276 |
| SiamMask + ASR | 0.379 | 0.604 | 0.280 |
| SiamMask + MD + ASR | 0.382 | 0.610 | 0.268 |
| SiamMask + MP | 0.384 | 0.610 | 0.262 |
| SiamMask + MD + MP | 0.394 | 0.611 | 0.234 |
| Ours | 0.397 | 0.612 | 0.220 |

TABLE VI

ABLATION STUDIES FOR MOTION PREDICTION AND ADAPTIVE SEARCH REGION ON VOT-2018 DATASET.

missing. For example, as for the car scenario in Fig. 6, when the camera shakes, the center of the search region of SiamRPN++ will shift to the left of the tracking car, and finally catches the truck. On the contrary, the center of our search region stays on the tracking car, since our model considers camera motion. This stability comes from the decoupling of camera motion. Another example is Bolt, the second row in Fig. 6. When Bolt accelerates, SiamRPN++ will be easily distracted by other runners, but our model does not fail because it considers the speed of Bolt. This stability comes from object velocity estimation. These unique features contribute to the performance of our method under large camera motion, fast object motion, and occlusion. In short, by predicting object position accurately, our model can focus on a more targeted search region and thus achieve better detection and segmentation performance.

*C. Evaluation for VOS*

Although this paper focuses on VOT, we also test our method on VOS.

*a) Datasets and settings:* We also report the performance of our method on standard VOS datasets DAVIS-2016 [20] and DAVIS-2017 [21]. For both datasets, we use the official performance measures: the Jaccard index (J) to express region similarity and the F-measure (F) to express contour accuracy. We use SiamMask as our segmentation module and adopt the semi-supervised setup. We fit bounding boxes to object masks in the first frame and use these bounding boxes to initialize our tracker.

*b) Results on DAVIS-2016 and DAVIS-2017:* Table V presents the comparison of VOS results using SiamMask and our proposed motion prediction model on Davis-2016 and Davis-2017 datasets. The effectiveness of our approach is limited on Davis-2016 and Davis-2017 datasets. The main reason is that DAVIS datasets have less camera motion or fast object motion. However, segmentation can still benefit from more accurately cropped search region. *e.g.,*, The dog in the third frame of the "Dogs-Jump" video is segmented more completely through motion prediction. However, SiamMask

misses the tail of the same dog during segmentation. Another example is the person in the fourth frame of the "Soap-Box" video. Our method separates this person from the soapbox. However, SiamMask mixes its segmentation with the surrounding pixels. Further, SiamMask fails to distinguish the person mask from the drum of the soapbox because the drum occupies the previous position of the person, which can not be handled without motion assumption. Though our pre-tracking procedure, our method can separate specific instance from its neighboring instance and thus get a more accurate segmentation. We show that our proposed method does a better job at segmenting under crowded scenarios. For more qualitative results, please refer to Fig. 7.

*D. Ablation studies*

Table VI compares the contribution of each module in our pipeline. Based on VOT-2018 dataset, we evaluate motion decouple (MD), motion prediction (MP), and adaptive search region (ASR) with the baseline approach (SiamMask). It can be observed from Table VI that the motion decouple and motion prediction play important roles in our method. The adaptive search region module only contributes 0.02 EAO improvement and using adaptive search region only even decrease the performance. This is because without motion decouple, the motion velocity might be very noisy. However, as we can see from Table VI, both of motion prediction and adaptive search region have the potential to improve accuracy with correct motion decouple.

## V. CONCLUSION

In conclusion, we show that motion prediction can still play an important role in visual object tracking. We propose a method that combines motion prediction with off-the-shelf appearance-based trackers. Although our baseline system is a straight forward combination of standard methods, we obtain the state-of-the-art results on VOT. We also show consistent improvements on VOS. We hope our work can inspire more
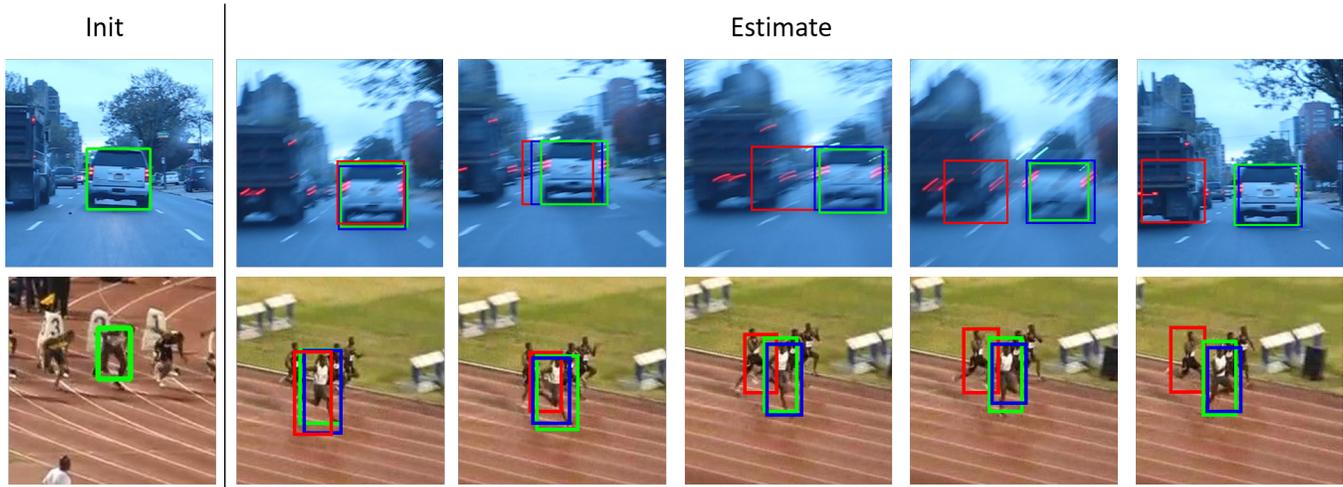
Fig. 6. Qualitative result of SiamRPN++ and our method : green box is the ground truth, red box is the bounding box from SiamRPN++, and blue box is our method.
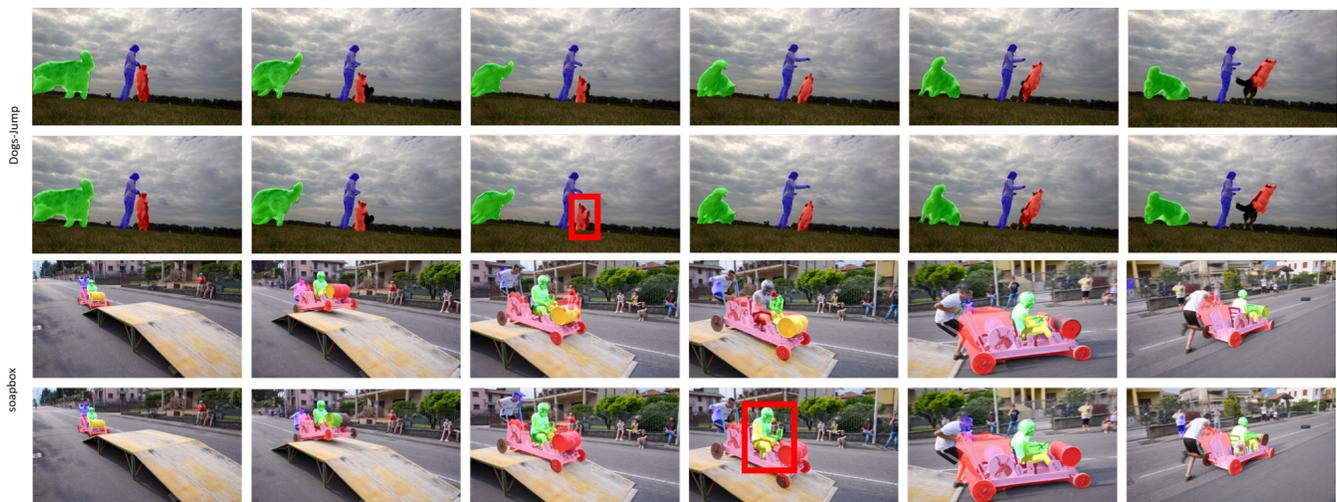


Fig. 7. Qualitative result of SiamMask and our method on DAVIS: First row and third row are the results from SiamMask. Second row and fourth row are the results from same videos using our method. (better view with color)

studies in considering the relationship between appearance and motion information in modern trackers.

## REFERENCES

[1] Q. Wang, L. Zhang, L. Bertinetto, W. Hu, and P. H. Torr, "Fast online object tracking and segmentation: A unifying approach," *arXiv preprint arXiv:1812.05050*, 2018.

[2] K.-H. Lee and J.-N. Hwang, "On-road pedestrian tracking across multiple driving recorders," *IEEE Transactions on Multimedia*, vol. 17, no. 9, pp. 1429–1438, 2015.

[3] S. Tang, M. Andriluka, B. Andres, and B. Schiele, "Multiple people tracking by lifted multicut and person reidentification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3539–3548.

[4] M. Mueller, N. Smith, and B. Ghanem, "A benchmark and simulator for uav tracking," in *European conference on computer vision*. Springer, 2016, pp. 445–461.

[5] Y. Kanayama, Y. Kimura, F. Miyazaki, and T. Noguchi, "A stable tracking control method for an autonomous mobile robot," in *Proceedings., IEEE International Conference on Robotics and Automation*, May 1990, pp. 384–389 vol.1.

[6] R. A. Singer, "Estimating optimal tracking filter performance for manned maneuvering targets," *IEEE Transactions on Aerospace and Electronic Systems*, vol. AES-6, no. 4, pp. 473–483, July 1970.

[7] D. Koller, K. Daniilidis, T. Thórhallson, and H. H. Nagel, "Model-based object tracking in traffic scenes," in *Computer Vision — ECCV'92*, G. Sandini, Ed. Berlin, Heidelberg: Springer Berlin Heidelberg, 1992, pp. 437–452.

[8] B. F. La Scala and R. R. Bitmead, "Design of an extended kalman filter frequency tracker," *IEEE Transactions on Signal Processing*, vol. 44, no. 3, pp. 739–742, March 1996.

[9] E. Maggio and A. Cavallaro, "Hybrid particle filter and mean shift tracker with adaptive transition model," in *Proceedings. (ICASSP '05). IEEE International Conference on Acoustics, Speech, and Signal Processing, 2005.*, vol. 2, March 2005, pp. ii/221–ii/224 Vol. 2.

[10] M. Isard and J. MacCormick, "Bramble: a bayesian multiple-blob tracker," in *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, vol. 2, July 2001, pp. 34–41 vol.2.

[11] S. Pellegrini, A. Ess, K. Schindler, and L. van Gool, "You'll never walk alone: Modeling social behavior for multi-target tracking," in *2009 IEEE 12th International Conference on Computer Vision*, Sep. 2009, pp. 261–268.

[12] A. Johansson, D. Helbing, and P. K. Shukla, "Specification of the

social force pedestrian model by evolutionary adjustment to video tracking data," *Advances in complex systems*, vol. 10, no. supp02, pp. 271–288, 2007.

[13] M. Kristan, J. Matas, A. Leonardis, M. Felsberg, G. Fernández, and et al., "The visual object tracking vot2016 challenge results," in *Proceedings of the European Conference on Computer Vision Workshop*. Springer International Publishing, 2016, pp. 777–823.

[14] M. Kristan, A. Leonardis, J. Matas, M. Felsberg, R. Pfugfelder, L. C. Zajc, T. Vojir, G. Bhat, A. Lukezic, A. Eldesokey, G. Fernandez, and et al., "The sixth visual object tracking vot2018 challenge results," 2018.

[15] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[16] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, 2015, pp. 91–99.

[17] B. Li, J. Yan, W. Wu, Z. Zhu, and X. Hu, "High performance visual tracking with siamese region proposal network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8971–8980.

[18] B. Li, W. Wu, Q. Wang, F. Zhang, J. Xing, and J. Yan, "Siamrpn++: Evolution of siamese visual tracking with very deep networks," *arXiv preprint arXiv:1812.11703*, 2018.

[19] Z. Zhu, Q. Wang, B. Li, W. Wu, J. Yan, and W. Hu, "Distractor-aware siamese networks for visual object tracking," in *European Conference on Computer Vision*. Springer, 2018, pp. 103–119.

[20] 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 Conference on Computer Vision and Pattern Recognition*, 2016, pp. 724–732.

[21] J. Pont-Tuset, F. Perazzi, S. Caelles, P. Arbeláez, A. Sorkine-Hornung, and L. Van Gool, "The 2017 davis challenge on video object segmentation," *arXiv preprint arXiv:1704.00675*, 2017.

[22] D. S. Bolme, J. R. Beveridge, B. A. Draper, and Y. M. Lui, "Visual object tracking using adaptive correlation filters," in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. IEEE, 2010, pp. 2544–2550.

[23] C. Ma, X. Yang, C. Zhang, and M.-H. Yang, "Long-term correlation tracking," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 5388–5396.

[24] Y. Li and J. Zhu, "A scale adaptive kernel correlation filter tracker with feature integration," in *European conference on computer vision*. Springer, 2014, pp. 254–265.

[25] H. K. Galoogahi, A. Fagg, and S. Lucey, "Learning background-aware correlation filters for visual tracking." in *ICCV*, 2017, pp. 1144–1152.

[26] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista, "High-speed tracking with kernelized correlation filters," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 3, pp. 583–596, 2015.

[27] M. Danelljan, F. Shahbaz Khan, M. Felsberg, and J. Van de Weijer, "Adaptive color attributes for real-time visual tracking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1090–1097.

[28] H. Kiani Galoogahi, T. Sim, and S. Lucey, "Multi-channel correlation filters," in *Proceedings of the IEEE international conference on computer vision*, 2013, pp. 3072–3079.

[29] M. Danelljan, G. Häger, F. Khan, and M. Felsberg, "Accurate scale estimation for robust visual tracking," in *British Machine Vision Conference, Nottingham, September 1-5, 2014*. BMVA Press, 2014.

[30] M. Danelljan, G. Häger, F. S. Khan, and M. Felsberg, "Discriminative scale space tracking," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 8, pp. 1561–1575, 2017.

[31] M. Danelljan, G. Hager, F. Shahbaz Khan, and M. Felsberg, "Learning spatially regularized correlation filters for visual tracking," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 4310–4318.

[32] H. Kiani Galoogahi, T. Sim, and S. Lucey, "Correlation filters with limited boundaries," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 4630–4638.

[33] M. Danelljan, A. Robinson, F. S. Khan, and M. Felsberg, "Beyond correlation filters: Learning continuous convolution operators for visual tracking," in *European Conference on Computer Vision*. Springer, 2016, pp. 472–488.

[34] G. Koch, R. Zemel, and R. Salakhutdinov, "Siamese neural networks for one-shot image recognition," in *ICML Deep Learning Workshop*, vol. 2, 2015.

[35] R. Tao, E. Gavves, and A. W. Smeulders, "Siamese instance search for tracking," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1420–1429.

[36] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. Torr, "Fully-convolutional siamese networks for object tracking," in *European conference on computer vision*. Springer, 2016, pp. 850–865.

[37] J. Valmadre, L. Bertinetto, J. Henriques, A. Vedaldi, and P. H. Torr, "End-to-end representation learning for correlation filter based tracking," in *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*. IEEE, 2017, pp. 5000–5008.

[38] D. Held, S. Thrun, and S. Savarese, "Learning to track at 100 fps with deep regression networks," in *European Conference on Computer Vision*. Springer, 2016, pp. 749–765.

[39] A. He, C. Luo, X. Tian, and W. Zeng, "A twofold siamese network for real-time object tracking," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[40] ——, "Towards a better match in siamese network based visual object tracker," in *European Conference on Computer Vision*. Springer, 2018, pp. 132–147.

[41] P. Luc, N. Neverova, C. Couprie, J. Verbeek, and Y. LeCun, "Predicting deeper into the future of semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 648–657.

[42] J. Walker, K. Marino, A. Gupta, and M. Hebert, "The pose knows: Video forecasting by generating pose futures," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3332–3341.

[43] X. Jin, X. Li, H. Xiao, X. Shen, Z. Lin, J. Yang, Y. Chen, J. Dong, L. Liu, Z. Jie *et al.*, "Video scene parsing with predictive feature learning," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 5580–5588.

[44] D. Nilsson and C. Sminchisescu, "Semantic video segmentation by gated recurrent flow propagation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6819–6828.

[45] M. Jaderberg, K. Simonyan, A. Zisserman *et al.*, "Spatial transformer networks," in *Advances in neural information processing systems*, 2015, pp. 2017–2025.

[46] H. Wang and C. Schmid, "Action recognition with improved trajectories," in *2013 IEEE International Conference on Computer Vision*, Dec 2013, pp. 3551–3558.

[47] G. Willems, T. Tuytelaars, and L. Van Gool, "An efficient dense and scale-invariant spatio-temporal interest point detector," in *European conference on computer vision*. Springer, 2008, pp. 650–663.

[48] S. Gauglitz, T. Höllerer, and M. Turk, "Evaluation of interest point detectors and feature descriptors for visual tracking," *International journal of computer vision*, vol. 94, no. 3, p. 335, 2011.

[49] D. Sun, X. Yang, M.-Y. Liu, and J. Kautz, "Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8934–8943.

[50] M. A. Fischler and R. C. Bolles, "Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.

[51] P. O. Pinheiro, R. Collobert, and P. Dollár, "Learning to segment object candidates," in *Advances in Neural Information Processing Systems*, 2015, pp. 1990–1998.

[52] P. O. Pinheiro, T.-Y. Lin, R. Collobert, and P. Dollár, "Learning to refine object segments," in *European Conference on Computer Vision*. Springer, 2016, pp. 75–91.

[53] D. Ramanan, D. A. Forsyth, and A. Zisserman, "Tracking people by learning their appearance," *IEEE transactions on pattern analysis and machine intelligence*, vol. 29, no. 1, pp. 65–81, 2007.

[54] J. Wang and Y. He, "Physics-aware 3d mesh synthesis," in *2019 International Conference on 3D Vision (3DV)*, 2019, pp. 502–512.

[55] J. Wang, Z. Fang, and H. Zhao, "Alignnet: A unifying approach to audio-visual alignment," in *The IEEE Winter Conference on Applications of Computer Vision (WACV)*, March 2020.