

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

## Prediction-Tracking-Segmentation

Anonymous ICCV submission

Paper ID 4319

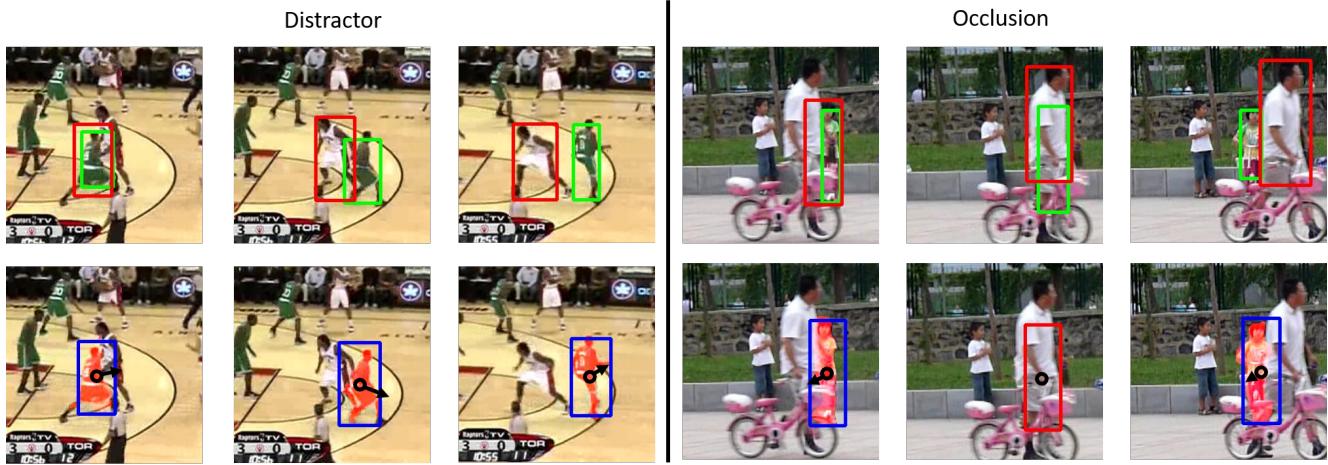


Figure 1: We propose a prediction driven method for tracking and segmentation in videos. The first row shows the results of state-of-the-art tracker SiamMask [56] (red) and ground truth (green). The second row shows the motion prediction (black arrow) and tracking results (with blue bounding box and red segmentation mask) of our Prediction-Tracking-Segmentation (PTS) model. Our method improves the robustness against distractors and occlusions. (better view with color)

## Abstract

We introduce a prediction driven method for visual tracking and segmentation in videos. Instead of solely relying on matching with appearance cues for tracking, we build a predictive model which provides guidance on finding more accurate tracking regions efficiently. With the proposed prediction mechanism, we improve the model robustness against distractions and occlusions during tracking. We demonstrate significant improvements over state-of-the-art methods not only on visual tracking tasks (VOT2016 and VOT2018) but also on video segmentation datasets (DAVIS2016 and DAVIS2017).

## 1. Introduction

Human can track, segment or interact with fast moving objects with surprising accuracy [46], even in the cases the objects are under deformations, occlusions and illumination changes [58]. What is the key component in human percep-

tion to make this happen?

In fact, tracking and segmenting moving object appears at a very early stage of human perception. Even 4-month-old infants can track moving objects with his or her eyes and reaching for them. [53]. A professional athlete can even interact with objects under very fast speed. (e.g. a baseball player can hit a 100-mph baseball). To do so, human brain has to overcome its own delays in neuronal transmission through prediction [38, 19, 6] and saves times for processing what we see by using only local information [51, 5]. Besides prediction, researchers further point out humans have multiple temporal scales for tracking objects with different speeds [20].

Inspired by the observations from human perception, in this paper, we propose a two-stage tracking method driven by prediction. Given a tracking result in time  $t$ , we first predicts the approximate object location in the next frame (in time  $t + 1$ ) without seeing it. Based on the prediction results, we then further refine the localization as well as the segmentation results by using the appearance input in time  $t + 1$ . The refined tracking and segmentation results can

108 help us back to update a better prediction model, which will  
109 be applied again in the successive frames.  
110

111 Specifically, in the first stage of prediction, we use an ex-  
112 trapolation method to estimate the object position in the fu-  
113 ture frame (in time  $t+1$ ), and simulate the multiple temporal  
114 scales effect in human perception through adaptive search  
115 region in the same frame. With prediction driven track-  
116 ing and segmentation, we name our method Prediction-  
117 Tracking-Segmentation (PTS).  
118

119 Our approach offers several unique advantages. First,  
120 through prediction of object position, we free tracking from  
121 using only appearance information. As most tracking and  
122 segmentation methods can only discriminate foreground  
123 from the non-semantic background [59], the performance  
124 suffers significantly when the target object is surrounded by  
125 similar objects (know as distractors [59]). Prediction also  
126 improves the robustness of our model against occlusions.  
127 Note that occlusions largely prevent most appearance based  
128 methods from extracting useful information. We show in  
129 the experiments that our method improves the tracking per-  
130 formance by a large margin under both cases. We visualize  
131 part of the results in Fig. 1.  
132

133 Second, through the usage of adaptive search region  
134 around the predicted area, PTS significantly decreases the  
135 information required to process and thus has a large poten-  
136 tial to increase the inference speed. To achieve this, we pro-  
137 pose to focus on smaller local regions when objects have  
138 slower speeds and smaller sizes, and vice versa. This ap-  
139 proach allows better segmentation performance, less miss-  
140 ing and identity switching.  
141

142 We evaluate our framework all major tracking datasets:  
143 VOT2016 [30], VOT2018 [29]. We demonstrate that  
144 our framework achieves state-of-the-art performance, both  
145 qualitatively and quantitatively. We also show competi-  
146 tive results against semi-supervised VOS approaches on  
147 DAVIS2016 [41] and DAVIS2017 [44].  
148

149 To summarize, our main contributions are three-fold:  
150 First, inspired by visual cognition theory, we propose PTS  
151 to unify predict, tracking and segmentation in a single  
152 framework. Second, we propose an adaptive search region  
153 module to effectively process information. We indicate that  
154 our proposed achieves competitive performance on VOT  
155 and VOS datasets.  
156

## 2. Related Works

157 In this section, we briefly overview three research areas  
158 relative to our proposed method.  
159

160 **Video Object Tracking** In tracking community, signif-  
161 icant attention has been paid to discriminative corre-  
162 lation filters (DCF) based methods [3, 35, 33, 13]. These  
163 methods allow to discriminate between the template of  
164 an arbitrary target and its 2D translations at a very fast  
165 speed. MOSSE [3] is the pioneer work which proposes  
166 a fast correlation tracker by minimizing the squared er-  
167 ror. Performance of DCF-based trackers has then been  
168 notably improved through the using of multi-channel fea-  
169 tures [18, 11, 26], robust scale estimation [7, 8], reducing  
170 boundary effects [9, 27] and fusing multi-resolution fea-  
171 tures in the continuous spatial domain [10].  
172

173 Tracking through Siamese Network is also an important  
174 approach [28, 48, 2, 50]. Instead of learning a discriminative  
175 classifier online, the idea is to train a deep siamese simila-  
176 rity function offline on pairs of video frames. At test time,  
177 this function is used to search for the candidate most similar  
178 to the template given in the starting frame on a new video,  
179 once per frame. The pioneering work is SINT [48]. Simi-  
180 larly, GOTURN [17] used deep regression network to pre-  
181 dict the motion between successive frames. SiamFC [2]  
182 implemented a fully convolutional network to output the  
183 correlation response map with high values at target loca-  
184 tions, which set a basic form of modern Siamese frame-  
185 work. Many following works have been proposed to im-  
186 prove the accuracy while maintain fast inference speed by  
187 adding semantic branch [16], using region proposals [32],  
188 hard negative mining [59], ensembling [15], deeper back-  
189 bone [31] and high-fidelity object representations [56].  
190

191 With the assumption that objects are under minor dis-  
192 placement and size change in consecutive frames, most  
193 modern trackers, including all the ones mentioned above,  
194 use a steady search region, which is centered on the last es-  
195 timated position of the target with the same ratio. Despite it  
196 is very straightforward, this oversimplified prior often fails  
197 in occlusion, motion change, size change, camera motion,  
198 as it is evident in the examples of Figure 1. This motivated  
199 us to propose a tracker able to adaptively set the search re-  
200 gion.  
201

202 **Video Forecasting** The ability to predict and therefore to  
203 anticipate the future is an important attribute of intelligence.  
204 Many methods are developed to improve the temporal sta-  
205 bility of semantic video segmentation. Luc et al. [34] de-  
206 velop an autoregressive convolutional neural network that  
207 learns to iteratively generate multiple future frames. Simi-  
208 larly, Walker et al. [54] uses a VAE to model the possi-  
209 ble future movements of humans in the pose space. In-  
210 stead of generating future states directly, many methods  
211 attempt to propagate segmentation from preceding input  
212 frames [24, 39, 21].  
213

214 Unlike previous work, inspired by human perception,  
215 we extract a motion model for each object, and setup new  
216 search region for segmentation according to the motion  
217 model.  
218

**216 Video Object Segmentation** Video Object Segmentation (VOS) have been divided into three categories based  
**217** on the level of supervision required: unsupervised, semi-  
**218** supervised and supervised. We briefly review the VOS fo-  
**219** cusing on semi-supervised setting, which is usually formu-  
**220** lated as a temporal label propagation problem. In order to  
**221** exploit consistency between video frames, many methods  
**222** propagate the first segmentation mask through temporal ad-  
**223** jacent ones [1, 36, 49] or even entire video [22, 23]. Another  
**224** approach is to process video frames independently [40]  
**225** and usually heavily rely on fine-tuning [4], data augmen-  
**226** tation [25] and model adaption [52].

### 228 3. Method

231 To unify prediction, tracking and segmentation with  
 232 adaptive search region, our model consists of: (i) predic-  
 233 tion module: estimate object position and velocity in an  
 234 unseen frame (ii) tracking module: adaptively limit the  
 235 search region for further processing (iii) segmentation mod-  
 236 ule: a fully-convolutional Siamese framework to segment  
 237 foreground object from given search region. We show our  
 238 framework in Figure 2.

#### 239 3.1. Prediction Module

240 By definition, the larger search image  $x$  is a larger crop  
 241 centered on the last estimated position of the target, which  
 242 is adopted by most Siamese trackers [2, 32, 59, 56].

243 However, these trackers only consider appearance fea-  
 244 tures of current frame, and hardly benefit from motion in-  
 245 formation. This leads to great difficulty in distinguishing  
 246 between instances that look like the template, known as dis-  
 247 tractors [59] or under occlusion, fast motion and camera  
 248 motion. To solve this problem, our proposed a tracker takes  
 249 full advantage of the motion information.

250 Object motion in a given image is the superposition of  
 251 camera motion and object motion, the former is random  
 252 while the latter should satisfy Newton’s First Law [37]. We  
 253 first pick a reference frame ( $F_{r_k}$ ,  $r_k$  denotes  $k^{th}$  refer-  
 254 ence frame) every  $n$  frames and thus separate the long video into  
 255 several pieces of short n-frame videos.

256 Second, we adopts the method proposed by ARIT [55]  
 257 to decouple the camera motion and object motion. ARIT  
 258 assumes that pending detection frame ( $F_{r_k+t}$ ) and its refer-  
 259 ence frame ( $F_{r_k}$ ) are related by a homography ( $H_{r_k, r_k+t}$ ).  
 260 This assumption holds in most cases as the global motion  
 261 between neighbouring frames is usually small. To estimate  
 262 the homography, the first step is to find the correspondences  
 263 between two frames. As mentioned in ARIT, we combine  
 264 SURF features [57] and motion vectors from the optical  
 265 flow in order to generate sufficient and complementary can-  
 266 didate matches, which is shown to be robust [14, 55]. Here  
 267 we use PWCNet [47] for dense flow generation.

268 As a homography matrix contains 8 free variables, at  
 269 least 4 background points pairs should be used. We cal-  
 270 culate the least square solution of Eq. 1 and optimize it to  
 271 obtain robust solution through RANSAC [12], where  $p_{r_k}^{bp}$   
 272 and  $p_{r_k+t}^{bp}$  denotes random selected back ground matching  
 273 pairs in  $F_{r_k}$  and  $F_{r_k+t}$  using the above mentioned features.  
 274 Given the assumption that the background occupies area  
 275 more than half of the images, we partition matching points  
 276 between frames into 4 pieces, then one point is randomly  
 277 chosen inside each selected piece to improve the efficiency  
 278 of RANSAC algorithm.

$$H_{r_k, r_k+t} \times p_{r_k}^{bp} = p_{r_k+t}^{bp} \quad (1)$$

280 For simplicity, all following calculations are under ref-  
 281 erence coordinate and project back to new coming frame  
 282 without further noticing.

283 Fig.3 illustrates the working principle of decoupling  
 284 step. The origin video for Fig.3 is a handheld video with  
 285 trembling background. The motion of the pedestrians in the  
 286 origin video is highly unpredictable with huge background  
 287 uncertainties. However, by mapping the target frame to-  
 288 wards the reference frame, the movement for pedestrians  
 289 could be more predictable and continuous.

290 In order to find the most representative point of object  
 291 position, we calculate the "center of mass" of object seg-  
 292 mentation using Eq. 2

$$P = \text{average}(p^o) \quad (2)$$

293 Random noise from background motion estimation and  
 294 mask segmentation might be introduced to the object pos-  
 295 ition prediction, which could influence the accuracy of pre-  
 296 diction. In order to achieve a better estimation for object  
 297 states, we utilize Kalman Filter to provide accurate pos-  
 298 ition information based on the measurements from current  
 299 and former frames. As a classical tracking algorithm, the  
 300 Kalman filter estimates the position of object in two steps:  
 301 prediction and correction. In prediction step, it predicts the  
 302 target state based on the dynamic model (eq.3) and gener-  
 303 ates the search region for Siamese network to achieve object  
 304 segmentation. Therefore, the measurement for object pos-  
 305 ition in the next frame could be computed with eq.2. Then,  
 306 in the correction step, the position measurement would be  
 307 updated with higher certainty given the position measure-  
 308 ment from Siamese network, which benefits the accuracy of  
 309 predictions for future frames.

310 The dynamic model for object position update could be  
 311 formulated as:

$$\hat{x}_{t|t-1} = F_k \hat{x}_{t-1|t-1} + w_t \quad (3)$$

312 In eq.3,  $\hat{x}_{t|t-1}$  is the priori state estimation given obser-  
 313 vations up to time  $t-1$ , which is in the form of 4-dimension  
 314 vector ( $[x, y, dx, dy]$ ) with position information. It is worth

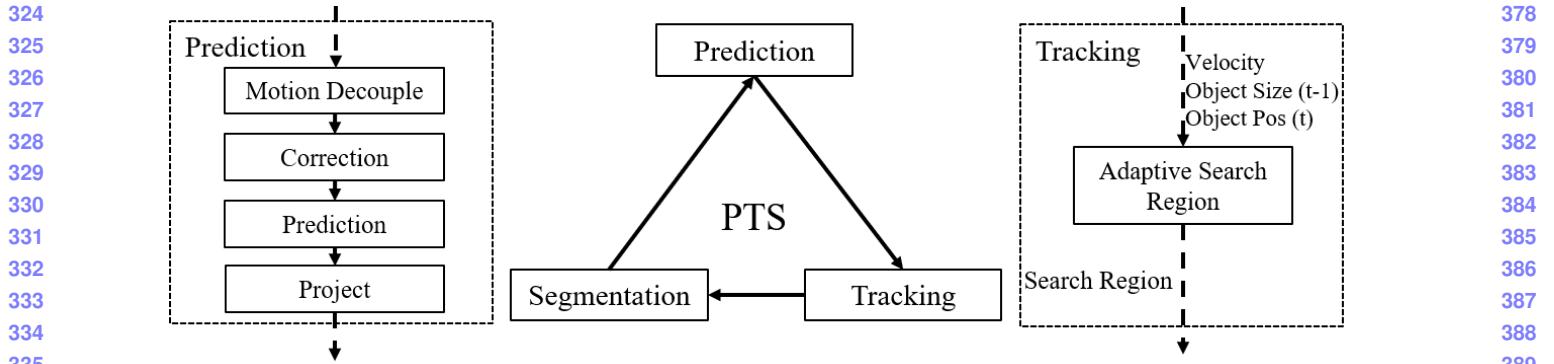


Figure 2: An overview of our method. Our method is composed of prediction part, tracking part and segmentation part.



Figure 3: One example for decoupling background motion and mapping object motion to reference frame (arrows illustrates the movement of object center)

to mention that the velocity terms ( $dx, dy$ ) are predicted by using extrapolation between the information from time  $t - 1$  and  $t$ .  $w_t$  is the random noise existing in the system. And  $F_t$  is the transition matrix from time  $t - 1$  to  $t$ .

After predicting the states, the Kalman filter uses measurements to correct its prediction during the correction steps using eq.4. In the equation,  $\hat{y}_t$  is the residuals between the prediction and measurement. And  $K_t$  is the optimal Kalman gain given from the predicted error covariance ( $P_{t|t-1}$ ), measurement matrix ( $H_t$ ) and measurement margin covariance ( $S_t$ ), as shown in eq.5. It is worth to mention that, as Kalman filter is a recursive algorithm, the predicted error covariance ( $P$ ) should be updated as well based on the estimation results.

$$\hat{x}_{t|t} = \hat{x}_{t|t-1} + K_t \hat{y}_t \quad (4)$$

$$K_t = P_{t|t-1} H_t^T S_t^{-1} \quad (5)$$

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 reference frame could be critical to the accuracy of 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 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.

### 3.2. Tracking Module

Inspired by human perception, we dynamically set up a new search region in the coming frame centered at the predicted object position. We project the estimated object center position back to the pending detection frame using Eq. 6.

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

Given the estimated position, we setup the search region  $S$  accordingly using the similar method as in [18]:

$$S = k \sqrt{(w + p)(h + p)} \quad (7)$$

$$k = 1 + 2 \times \text{sigmoid}(\|v\|_2 - T) \quad (8)$$

where  $p = \frac{w + h}{2}$ . To achieve the adaptive search region, the search region size  $S$  would be modified with respect to the predicted velocity using eq.8. In the equation,  $v$  is the velocity predicted by Kalman filter and  $T$  is the threshold for velocity. The search region is cropped center at  $P_{r_k+t}$  on the frame  $F_{r+k}$ , and then resized in  $255 \times 255$ .

To make the one-shot segmentation framework suitable for tracking task. We adopt the optimisation strategy used for the automatic bounding box generation proposed in

432 VOT2016 [30] as it offers the highest IOU and mAP as re-  
 433 ported in [56].

### 434 435 3.3. Segmentation Module

436 We adopt the SiamMask framework [56], which achieves  
 437 a good balance between the accuracy and speed. SiamMask  
 438 propose to use an offline trained fully-convolutional net-  
 439 work to simultaneously collect binary segmentation mask,  
 440 detection bounding box and objectness score. First, the  
 441 Siamese network compares an template image  $z$  ( $w \times h \times 3$ )  
 442 against a (larger) search image  $x$  to obtain a dense re-  
 443 sponse map  $g_\theta(z, x)$ . The two inputs are processed by the  
 444 same CNN  $f_\theta$ , yielding two feature maps that are cross-  
 445 correlated:

$$446 g_\theta(z, x) = f_\theta(z) * f_\theta(x) \quad (9)$$

447 Each spatial element of the response map  $g_\theta^n(z, x)$  represent  
 448 a similarity between the template image  $z$  and  $n^{th}$  candidate  
 449 window in  $x$ . Second, a three-branch head calculates binary  
 450 segmentation mask, detection bounding box and objectness  
 451 score, respectively. The mask branch predicts a ( $w \times h$ ) bi-  
 452 nary mask  $m^n$  from each spatial element  $g_\theta^n(z, x)$ . The box  
 453 branch regresses  $k$  bounding boxes from each spatial ele-  
 454 ment  $g_\theta^n(z, x)$ , where  $k$  is the number of anchors. And the  
 455 score branch estimates the corresponding objectness score.  
 456

$$457 m^n = h_\phi(g_\theta^n(z, x)) \quad (10)$$

$$458 b_{i=1,\dots,k}^n = h_\sigma(g_\theta^n(z, x)) \quad (11)$$

$$459 s_{i=1,\dots,k}^n = h_\varphi(g_\theta^n(z, x)) \quad (12)$$

460 A multi-task loss is used to optimise the whole framework.

$$461 L = \lambda_1 L_{mask} + \lambda_2 L_{box} + \lambda_3 L_{score} \quad (13)$$

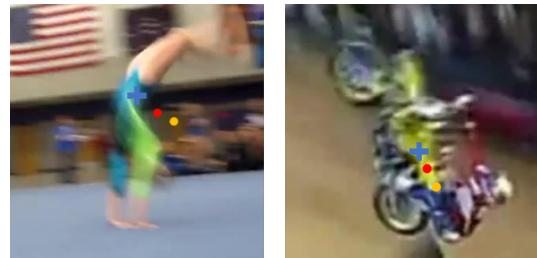
462 We refer readers to [42, 43] for understanding mask  
 463 branch and [45, 32] for understanding region proposal  
 464 branch.

## 465 466 4. Experiments

467 In this section, we evaluate our approach on three tasks:  
 468 object position prediction, visual object tracking (VOT2016  
 469 and VOT2018) and semi-supervised video object segmen-  
 470 tation (on DAVIS2016 and DAVIS2017).

### 471 472 4.1. Evaluation for prediction accuracy

473 **Datasets and settings** We adopt two widely used bench-  
 474 mark data set to evaluate the performance the motion pre-  
 475 diction module: VOT2016 [30] and VOT2018 [29], which  
 476 both of them are annotated with rotated bounding box. We  
 477 choose these two dataset because they contain complicated  
 478 camera motion and object motion, which is extremely chal-  
 479 lenging for motion prediction [29]. We use mean distance  
 480 from bounding box center as prediction error. We then com-  
 481 pare against the state-of-the-art tracker.



482 Figure 4: Object center predictions generated by SiamMask  
 483 and PTS (red for PTS, yellow for SiamMask and blue cross  
 484 for ground truth)

485 **Results on VOT2016 and VOT2018** Table.2 presents the  
 486 comparison of prediction results using SiamMask and PTS  
 487 based on VOT2016 and VOT2018 datasets. As it is shown  
 488 in the table, for both of these two datasets, PTS method  
 489 could dramatically reduce the prediction errors of the ob-  
 490 ject position. The mean square error for object position on  
 491 VOT2018 could be reduced to half from 16 pixels to 8 pix-  
 492 els. Meanwhile, Fig.4 shows when the object velocity is  
 493 high, PTS method could provide a prediction more accurate  
 494 compared with the SiamMask, which does not consider the  
 495 influence of object motion. The results prove that the decou-  
 496 pling strategy could reduce the background uncertainty and  
 497 the Kalman filter would provide relatively reliable predic-  
 498 tion for object position in the next frame. Higher accuracy  
 499 for object position prediction could benefit the generation  
 500 of search regions for object tracking and eventually improve  
 501 the performance of object segmentation.

### 502 503 4.2. Evaluation for visual object tracking

504 **Datasets and settings** Similarly, we adopt three widely  
 505 used benchmarks for the evaluation of the object tracking  
 506 task: VOT2016, VOT2018 and compare against the state-  
 507 of-the-art using official metric: Expected Average Overlap  
 508 (EAO). We use VOT2018 to conduct an experiment to fur-  
 509 ther discuss the performance under different conditions.

510 **Results on VOT2016 and VOT2018** Table 1 and Table  
 511 3 present comparisons of tracking performance between  
 512 PTS and other state of art models based on VOT2016 and  
 513 VOT2018 datasets. Our model significantly improves the  
 514 robustness, and achieves a EAO of 0.397.

515 The robustness comes from the unique feature of PTS  
 516 which can greatly help with large camera motion, fast object  
 517 motion and occlusion. Comparing with SiamMask, which  
 518 uses the position generated from last mask as the center of  
 519 search region for current frame, we use adaptive search re-  
 520 gion module and motion prediction module to predict the  
 521 center of search region, so the center of the search region  
 522 will be more likely to locate on the tracking object. For

		DaSiamRPN	SA_Siam_R	CPT	DeepSTRCF	DRT	RCO	UPDT	SiamMask	SiamRPN	MFT	LADCF	Ours	
540	EAO	0.326	0.337	0.339	0.345	0.356	0.376	0.378	0.380	0.383	0.385	0.389	<b>0.397</b>	594
541	Accuracy	0.569	0.566	0.506	0.523	0.519	0.507	0.536	0.609	0.586	0.505	0.503	0.612	595
542	Robustness	0.337	0.258	0.239	0.215	0.201	0.155	0.184	0.276	0.276	0.140	0.159	0.220	596
543														597

Table 1: Comparison with the state-of-the-art under EAO, Accuracy, and Robustness on the VOT-2018 benchmark.

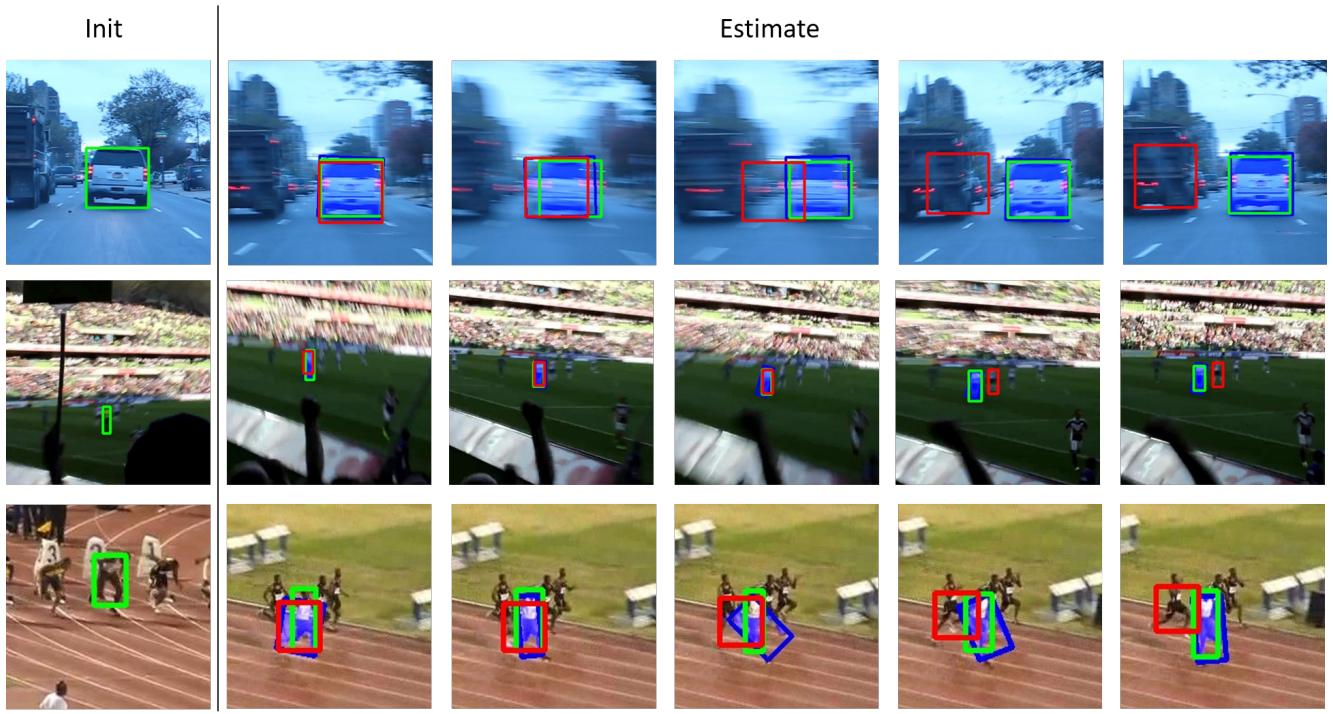


Figure 5: Qualitative result of our method : green box is the ground truth, red box is the bounding box from SiamMask, and blue box is our bounding box for the mask.

VOT2016		VOT2018	
Tracker	MSE	Tracker	MSE
SiamMask	16.281	SiamMask	14.593
PST(ours)	8.198	PST(ours)	7.544

Table 2: Mean square error for prediction on VOT2016 and VOT2018

Trackers	VOT2016		
	A	R	EAO
SiamMaks	0.639	0.214	0.433
PST (Ours)	0.642	0.144	0.471

Table 3: Comparison with SiamMask on VOT 2016

example, as for the car scenario in figure 5, when the camera shakes, the center of search region of SiamMask will

shift to the left of the tracking car, and finally catches the truck. On the contrary, since our model considers camera motion, the center of our search region stays on the tracking car. Another example is Bolt, the third row in figure 5. When Bolt accelerates, SiamMask will be easily distracted by other runners, but our PTS model won't fail because it considers the speed of Bolt.

### 4.3. Evaluation for video object segmentation

**Datasets and settings** We report the performance of PWT on standard VOS datasets DAVIS2016 [41] and DAVIS2017 [44]. 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.

**Results on DAVIS2016 and DAVIS2017** Table. 4 presents the comparison of vos results using SiamMask and PTS based on DAVIS2016 and DAVIS2017 datasets.

	J	F
DAVIS2016	0.732	0.692
DAVIS2017	0.554	0.604

Table 4: J and F Results on DAVIS2016 and DAVIS2017

	EAO	Accuracy	Robustness
SiamMask	0.380	0.609	0.276
SiamMask + Tracking	0.382	0.610	0.268
SiamMask + Prediction	0.394	0.611	0.234
PST	0.397	0.612	0.220

Table 5: Ablation studies for Tracking and Prediction modules on VOT2018 dataset.

The affect of our approach is limited on DAVIS2016 and 2017 dataset. The main reason is that DAVIS dataset has less camera motion or fast object motion, which are the major gain from our method.

#### 4.4. Ablation studies

Table 5 compares the influence of each modules in our model. Based on VOT2018 dataset, we evaluate the influence of tracking and prediction module and compare their performance with the baseline approach (SiamMask) and PST. It can be observed from Table 5 that the prediction module which uses Kalman filter to update the position of objects plays an important role in PST that most EAO improvements seems to be introduced by prediction module. And for the tracking module, which is the adaptive search region update module, the influence is a little bit limited with only 0.02 EAO increase. But as we can see from Table 5, both of this two module has the potential to improve the accuracy.

### 5. Conclusion

In conclusion, we introduce a prediction driven method for visual tracking and segmentation in videos. Instead of solely relying on matching with appearance cues for tracking, we build a predictive model which provides guidance on finding more accurate tracking regions efficiently. With the pro-posed prediction mechanism, we improve the model robustness against distractions and occlusions during tracking. We demonstrate significant improvements over state-of-the-art methods not only on visual tracking tasks (VOT2016 and VOT2018) but also on video segmentation datasets (DAVIS2016 and DAVIS2017).

### References

- [1] L. Bao, B. Wu, and W. Liu. Cnn in mrf: Video object segmentation via inference in a cnn-based higher-order spatio-temporal mrf. In *Proceedings of the IEEE Conference*

- on Computer Vision and Pattern Recognition
- 702
- 5977–5986, 2018. 3
- 703
- [2] 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*, pages 850–865. Springer, 2016. 2, 3
- 704
- [3] 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*, pages 2544–2550. IEEE, 2010. 2
- 705
- [4] S. Caelles, K.-K. Maninis, J. Pont-Tuset, L. Leal-Taixé, D. Cremers, and L. Van Gool. One-shot video object segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 221–230, 2017. 3
- 706
- [5] C. R. Cassanello, A. T. Nihalani, and V. P. Ferrera. Neuronal responses to moving targets in monkey frontal eye fields. *Journal of neurophysiology*, 2008. 1
- 707
- [6] P. Cavanagh and S. Anstis. The flash grab effect. *Vision Research*, 91:8–20, 2013. 1
- 708
- [7] 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. 2
- 709
- [8] M. Danelljan, G. Häger, F. S. Khan, and M. Felsberg. Discriminative scale space tracking. *IEEE transactions on pattern analysis and machine intelligence*, 39(8):1561–1575, 2017. 2
- 710
- [9] 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*, pages 4310–4318, 2015. 2
- 711
- [10] 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*, pages 472–488. Springer, 2016. 2
- 712
- [11] 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*, pages 1090–1097, 2014. 2
- 713
- [12] 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*, 24(6):381–395, 1981. 3
- 714
- [13] H. K. Galoogahi, A. Fagg, and S. Lucey. Learning background-aware correlation filters for visual tracking. In *ICCV*, pages 1144–1152, 2017. 2
- 715
- [14] S. Gauglitz, T. Höllerer, and M. Turk. Evaluation of interest point detectors and feature descriptors for visual tracking. *International journal of computer vision*, 94(3):335, 2011. 3
- 716
- [15] A. He, C. Luo, X. Tian, and W. Zeng. Towards a better match in siamese network based visual object tracker. In *European Conference on Computer Vision*, pages 132–147. Springer, 2018. 2
- 717
- [16] 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. 2
- 718
- [17] C. R. Cassanello, A. T. Nihalani, and V. P. Ferrera. Neuronal responses to moving targets in monkey frontal eye fields. *Journal of neurophysiology*, 2008. 1
- 719
- [18] P. Cavanagh and S. Anstis. The flash grab effect. *Vision Research*, 91:8–20, 2013. 1
- 720
- [19] 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. 2
- 721
- [20] M. Danelljan, G. Häger, F. S. Khan, and M. Felsberg. Discriminative scale space tracking. *IEEE transactions on pattern analysis and machine intelligence*, 39(8):1561–1575, 2017. 2
- 722
- [21] 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*, pages 4310–4318, 2015. 2
- 723
- [22] 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*, pages 472–488. Springer, 2016. 2
- 724
- [23] 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*, pages 1090–1097, 2014. 2
- 725
- [24] 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*, 24(6):381–395, 1981. 3
- 726
- [25] H. K. Galoogahi, A. Fagg, and S. Lucey. Learning background-aware correlation filters for visual tracking. In *ICCV*, pages 1144–1152, 2017. 2
- 727
- [26] S. Gauglitz, T. Höllerer, and M. Turk. Evaluation of interest point detectors and feature descriptors for visual tracking. *International journal of computer vision*, 94(3):335, 2011. 3
- 728
- [27] A. He, C. Luo, X. Tian, and W. Zeng. Towards a better match in siamese network based visual object tracker. In *European Conference on Computer Vision*, pages 132–147. Springer, 2018. 2
- 729
- [28] 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. 2
- 730
- [29] C. R. Cassanello, A. T. Nihalani, and V. P. Ferrera. Neuronal responses to moving targets in monkey frontal eye fields. *Journal of neurophysiology*, 2008. 1
- 731
- [30] P. Cavanagh and S. Anstis. The flash grab effect. *Vision Research*, 91:8–20, 2013. 1
- 732
- [31] 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. 2
- 733
- [32] M. Danelljan, G. Häger, F. S. Khan, and M. Felsberg. Discriminative scale space tracking. *IEEE transactions on pattern analysis and machine intelligence*, 39(8):1561–1575, 2017. 2
- 734
- [33] 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*, pages 4310–4318, 2015. 2
- 735
- [34] 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*, pages 472–488. Springer, 2016. 2
- 736
- [35] 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*, pages 1090–1097, 2014. 2
- 737
- [36] 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*, 24(6):381–395, 1981. 3
- 738
- [37] H. K. Galoogahi, A. Fagg, and S. Lucey. Learning background-aware correlation filters for visual tracking. In *ICCV*, pages 1144–1152, 2017. 2
- 739
- [38] S. Gauglitz, T. Höllerer, and M. Turk. Evaluation of interest point detectors and feature descriptors for visual tracking. *International journal of computer vision*, 94(3):335, 2011. 3
- 740
- [39] A. He, C. Luo, X. Tian, and W. Zeng. Towards a better match in siamese network based visual object tracker. In *European Conference on Computer Vision*, pages 132–147. Springer, 2018. 2
- 741
- [40] 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. 2
- 742
- [41] C. R. Cassanello, A. T. Nihalani, and V. P. Ferrera. Neuronal responses to moving targets in monkey frontal eye fields. *Journal of neurophysiology*, 2008. 1
- 743
- [42] P. Cavanagh and S. Anstis. The flash grab effect. *Vision Research*, 91:8–20, 2013. 1
- 744
- [43] 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. 2
- 745
- [44] M. Danelljan, G. Häger, F. S. Khan, and M. Felsberg. Discriminative scale space tracking. *IEEE transactions on pattern analysis and machine intelligence*, 39(8):1561–1575, 2017. 2
- 746
- [45] 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*, pages 4310–4318, 2015. 2
- 747
- [46] 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*, pages 472–488. Springer, 2016. 2
- 748
- [47] 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*, pages 1090–1097, 2014. 2
- 749
- [48] 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*, 24(6):381–395, 1981. 3
- 750
- [49] H. K. Galoogahi, A. Fagg, and S. Lucey. Learning background-aware correlation filters for visual tracking. In *ICCV*, pages 1144–1152, 2017. 2
- 751
- [50] S. Gauglitz, T. Höllerer, and M. Turk. Evaluation of interest point detectors and feature descriptors for visual tracking. *International journal of computer vision*, 94(3):335, 2011. 3
- 752
- [51] A. He, C. Luo, X. Tian, and W. Zeng. Towards a better match in siamese network based visual object tracker. In *European Conference on Computer Vision*, pages 132–147. Springer, 2018. 2
- 753
- [52] 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. 2
- 754
- [53] C. R. Cassanello, A. T. Nihalani, and V. P. Ferrera. Neuronal responses to moving targets in monkey frontal eye fields. *Journal of neurophysiology*, 2008. 1
- 755

- 756 [17] D. Held, S. Thrun, and S. Savarese. Learning to track at 100  
757 fps with deep regression networks. In *European Conference  
758 on Computer Vision*, pages 749–765. Springer, 2016. 2
- 759 [18] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista. High-  
760 speed tracking with kernelized correlation filters. *IEEE  
761 Transactions on Pattern Analysis and Machine Intelligence*,  
762 37(3):583–596, 2015. 2, 4
- 763 [19] H. Hogendoorn and A. N. Burkitt. Predictive coding of vi-  
764 sual object position ahead of moving objects revealed by  
765 time-resolved eeg decoding. *Neuroimage*, 171:55–61, 2018.  
766 1
- 767 [20] A. O. Holcombe. Seeing slow and seeing fast: two limits  
768 on perception. *Trends in cognitive sciences*, 13(5):216–221,  
769 2009. 1
- 770 [21] M. Jaderberg, K. Simonyan, A. Zisserman, et al. Spatial  
771 transformer networks. In *Advances in neural information  
772 processing systems*, pages 2017–2025, 2015. 2
- 773 [22] V. Jampani, R. Gaddé, and P. V. Gehler. Video propagation  
774 networks. In *Proceedings of the IEEE Conference on Com-  
775 puter Vision and Pattern Recognition*, pages 451–461, 2017.  
776 3
- 777 [23] W.-D. Jang and C.-S. Kim. Online video object segmen-  
778 tation via convolutional trident network. In *Proceedings of the  
779 IEEE Conference on Computer Vision and Pattern Recog-  
780 nition*, pages 5849–5858, 2017. 3
- 781 [24] X. Jin, X. Li, H. Xiao, X. Shen, Z. Lin, J. Yang, Y. Chen,  
782 J. Dong, L. Liu, Z. Jie, et al. Video scene parsing with pre-  
783 dictive feature learning. In *Proceedings of the IEEE Interna-  
784 tional Conference on Computer Vision*, pages 5580–5588,  
785 2017. 2
- 786 [25] A. Khoreva, R. Benenson, E. Ilg, T. Brox, and B. Schiele.  
787 Lucid data dreaming for object tracking. In *The DAVIS Chal-  
788 lenge on Video Object Segmentation*, 2017. 3
- 789 [26] H. Kiani Galoogahi, T. Sim, and S. Lucey. Multi-channel  
790 correlation filters. In *Proceedings of the IEEE international  
791 conference on computer vision*, pages 3072–3079, 2013. 2
- 792 [27] H. Kiani Galoogahi, T. Sim, and S. Lucey. Correlation filters  
793 with limited boundaries. In *Proceedings of the IEEE Confer-  
794 ence on Computer Vision and Pattern Recognition*, pages  
795 4630–4638, 2015. 2
- 796 [28] G. Koch, R. Zemel, and R. Salakhutdinov. Siamese neu-  
797 ral networks for one-shot image recognition. In *ICML Deep  
798 Learning Workshop*, volume 2, 2015. 2
- 799 [29] M. Kristan, A. Leonardis, J. Matas, M. Felsberg,  
800 R. Pfugfelder, L. C. Zajc, T. Vojir, G. Bhat, A. Lukezic,  
801 A. Elde索key, G. Fernandez, and et al. The sixth visual  
802 object tracking vot2018 challenge results, 2018. 2, 5
- 803 [30] M. Kristan, J. Matas, A. Leonardis, M. Felsberg,  
804 G. Fernández, and et al. The visual object tracking vot2016  
805 challenge results. In *Proceedings of the European Con-  
806 ference on Computer Vision Workshop*, pages 777–823.  
807 Springer International Publishing, 2016. 2, 5
- 808 [31] B. Li, W. Wu, Q. Wang, F. Zhang, J. Xing, and J. Yan.  
809 Siamrpn++: Evolution of siamese visual tracking with very  
810 deep networks. *arXiv preprint arXiv:1812.11703*, 2018. 2
- 811 [32] B. Li, J. Yan, W. Wu, Z. Zhu, and X. Hu. High per-  
812 formance visual tracking with siamese region proposal network.  
813
- 814 [33] Y. Li and J. Zhu. A scale adaptive kernel correlation filter  
815 tracker with feature integration. In *European conference on  
816 computer vision*, pages 254–265. Springer, 2014. 2
- 817 [34] P. Luc, N. Neverova, C. Couprie, J. Verbeek, and Y. Le-  
818 Cun. Predicting deeper into the future of semantic segmen-  
819 tation. In *Proceedings of the IEEE International Conference  
820 on Computer Vision*, pages 648–657, 2017. 2
- 821 [35] C. Ma, X. Yang, C. Zhang, and M.-H. Yang. Long-term cor-  
822 relation tracking. In *Proceedings of the IEEE conference on  
823 computer vision and pattern recognition*, pages 5388–5396,  
824 2015. 2
- 825 [36] N. Märki, F. Perazzi, O. Wang, and A. Sorkine-Hornung.  
826 Bilateral space video segmentation. In *Proceedings of the  
827 IEEE Conference on Computer Vision and Pattern Recog-  
828 nition*, pages 743–751, 2016. 3
- 829 [37] I. Newton. *The Principia: mathematical principles of nat-  
830 ural philosophy*. Univ of California Press, 1999. 3
- 831 [38] R. Nijhawan. Motion extrapolation in catching. *Nature*,  
832 1994. 1
- 833 [39] D. Nilsson and C. Sminchisescu. Semantic video segmen-  
834 tation by gated recurrent flow propagation. In *Proceedings  
835 of the IEEE Conference on Computer Vision and Pattern  
836 Recognition*, pages 6819–6828, 2018. 2
- 837 [40] F. Perazzi, A. Khoreva, R. Benenson, B. Schiele, and  
838 A. Sorkine-Hornung. Learning video object segmentation  
839 from static images. In *Proceedings of the IEEE Conference  
840 on Computer Vision and Pattern Recognition*, pages 2663–  
841 2672, 2017. 3
- 842 [41] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool,  
843 M. Gross, and A. Sorkine-Hornung. A benchmark dataset  
844 and evaluation methodology for video object segmentation.  
845 In *Proceedings of the IEEE Conference on Computer Vision  
846 and Pattern Recognition*, pages 724–732, 2016. 2, 6
- 847 [42] P. O. Pinheiro, R. Collobert, and P. Dollár. Learning to seg-  
848 ment object candidates. In *Advances in Neural Information  
849 Processing Systems*, pages 1990–1998, 2015. 5
- 850 [43] P. O. Pinheiro, T.-Y. Lin, R. Collobert, and P. Dollár. Learn-  
851 ing to refine object segments. In *European Conference on  
852 Computer Vision*, pages 75–91. Springer, 2016. 5
- 853 [44] J. Pont-Tuset, F. Perazzi, S. Caelles, P. Arbeláez, A. Sorkine-  
854 Hornung, and L. Van Gool. The 2017 davis chal-  
855 lenge on video object segmentation. *arXiv preprint  
856 arXiv:1704.00675*, 2017. 2, 6
- 857 [45] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards  
858 real-time object detection with region proposal networks. In  
859 *Advances in neural information processing systems*, pages  
860 91–99, 2015. 5
- 861 [46] J. B. Smeets, E. Brenner, and M. H. de Lussanet. Visuomotor  
862 delays when hitting running spiders. In *EWEP 5-Advances  
863 in perception-action coupling*, pages 36–40. Éditions EDK,  
864 Paris, 1998. 1
- 865 [47] D. Sun, X. Yang, M.-Y. Liu, and J. Kautz. Pwc-net: Cnns  
866 for optical flow using pyramid, warping, and cost volume.  
867 In *Proceedings of the IEEE Conference on Computer Vision  
868 and Pattern Recognition*, pages 8934–8943, 2018. 3

- 864 [48] R. Tao, E. Gavves, and A. W. Smeulders. Siamese instance 918  
865 search for tracking. In *Proceedings of the IEEE conference 919*  
866 on computer vision and pattern recognition, pages 1420– 920  
867 1429, 2016. 2 921
- 868 [49] Y.-H. Tsai, M.-H. Yang, and M. J. Black. Video segmentation 922  
869 via object flow. In *Proceedings of the IEEE conference 923*  
870 on computer vision and pattern recognition, pages 3899– 924  
871 3908, 2016. 3 925
- 872 [50] J. Valmadre, L. Bertinetto, J. Henriques, A. Vedaldi, and 926  
873 P. H. Torr. End-to-end representation learning for correlation 927  
874 filter based tracking. In *Computer Vision and Pattern 928*  
875 Recognition (CVPR), 2017 IEEE Conference on, pages 5000–5008. 929  
876 IEEE, 2017. 2 930
- 877 [51] E. van Heusden, M. Rolfs, P. Cavanagh, and H. Hogendoorn. 931  
878 Motion extrapolation for eye movements predicts perceived 932  
879 motion-induced position shifts. *Journal of Neuroscience*, 933  
880 38(38):8243–8250, 2018. 1 934
- 881 [52] P. Voigtlaender and B. Leibe. Online adaptation of convolutional 935  
882 neural networks for video object segmentation. *arXiv 936*  
883 preprint arXiv:1706.09364, 2017. 3 937
- 884 [53] C. Von Hofsten. Eye-hand coordination in the newborn. *De- 938*  
885 velopmental psychology, 18(3):450, 1982. 1 939
- 886 [54] J. Walker, K. Marino, A. Gupta, and M. Hebert. The pose 940  
887 knows: Video forecasting by generating pose futures. In *Pro- 941*  
888 ceedings of the IEEE International Conference on Computer 942  
889 Vision, pages 3332–3341, 2017. 2 943
- 890 [55] H. Wang and C. Schmid. Action recognition with improved 944  
891 trajectories. In *2013 IEEE International Conference on 945*  
892 Computer Vision, pages 3551–3558, Dec 2013. 3 946
- 893 [56] Q. Wang, L. Zhang, L. Bertinetto, W. Hu, and P. H. Torr. 947  
894 Fast online object tracking and segmentation: A unifying 948  
895 approach. *arXiv preprint arXiv:1812.05050*, 2018. 1, 2, 3, 5 949
- 896 [57] G. Willems, T. Tuytelaars, and L. Van Gool. An efficient 950  
897 dense and scale-invariant spatio-temporal interest point 951  
898 detector. In *European conference on computer vision*, pages 952  
899 650–663. Springer, 2008. 3 953
- 900 [58] Y. Wu, J. Lim, and M.-H. Yang. Object tracking benchmark. 954  
901 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9):1834–1848, 2015. 1 955
- 902 [59] Z. Zhu, Q. Wang, B. Li, W. Wu, J. Yan, and W. Hu. 956  
903 Distractor-aware siamese networks for visual object tracking. 957  
904 In *European Conference on Computer Vision*, pages 958  
905 103–119. Springer, 2018. 2, 3 959
- 906 960
- 907 961
- 908 962
- 909 963
- 910 964
- 911 965
- 912 966
- 913 967
- 914 968
- 915 969
- 916 970
- 917 971