# Need for Speed: A Benchmark for Higher Frame Rate Object Tracking

Hamed Kiani Galoogahi[1]*, Ashton Fagg[2,1]*, Chen Huang[1], Deva Ramanan[1] and Simon Lucey[1,2]

[1]Robotics Institute
Carnegie Mellon University

[2] SAIVT Lab
Queensland University of Technology

## Abstract

*In this paper, we propose the first higher frame rate video dataset (called Need for Speed - NfS) and benchmark for visual object tracking. The dataset consists of 100 videos (380K frames) captured with now commonly available higher frame rate (240 FPS) cameras from real world scenarios. All frames are annotated with axis aligned bounding boxes and all sequences are manually labelled with nine visual attributes - such as occlusion, fast motion, background clutter, etc. Our benchmark provides an extensive evaluation of many recent and state-of-the-art trackers on higher frame rate sequences. We ranked each of these trackers according to their tracking accuracy and real-time performance. One of our surprising conclusions is that at higher frame rates, simple trackers such as correlation filters outperform complex methods based on deep networks. This suggests that for practical applications (such as in robotics or embedded vision), one needs to carefully trade-off bandwidth constraints associated with higher frame rate acquisition, computational costs of real-time analysis, and the required application accuracy. Our dataset and benchmark allows for the first time (to our knowledge) systematic exploration of such issues, and will be made available to allow for further research in this space.*

## 1. Introduction

Visual object tracking is a fundamental task in computer vision which has implications for a bevy of applications: surveillance, vehicle autonomy, video analysis, etc. The vision community has shown an increasing degree of interest in the problem - with recent methods becoming increasingly sophisticated and accurate [8, 13, 10, 1, 2]. However, most of these algorithms - and the datasets they have been evaluated upon [37, 36, 22] - have been aimed at the canonical approximate frame rate of 30 Frames Per Second (FPS). Consumer devices with cameras such as smart phones, tablets,
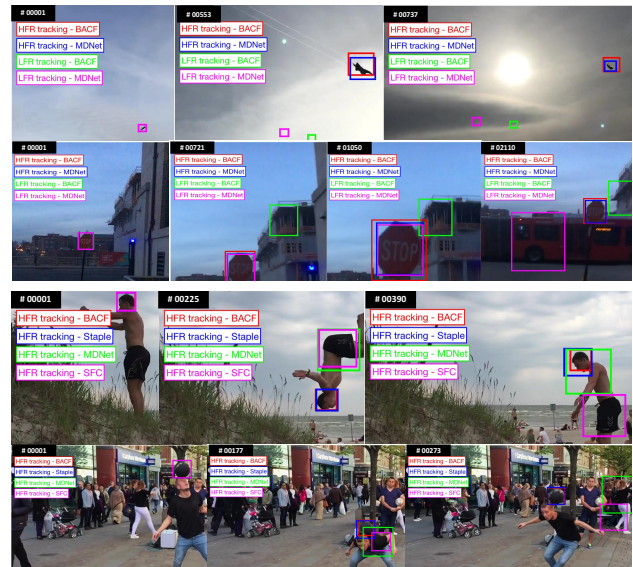


Figure 1. The effect of tracking higher frame rate videos. Top rows illustrate the robustness of tracking higher frame rate videos (240 FPS) versus lower frame rate videos (30 FPS) for a Correlation Filter (BACF with HOG) and a deep tracker (MDNet). Bottom rows show if higher frame rate videos are available, cheap CF trackers (Staple and BACF) can outperform complicated deep trackers (SFC and MDNet) on challenging situations such as fast motion, rotation, illumination and cluttered background. Predicted bounding boxes of these methods are shown by different colors. HFR and LFR refer to Higher and Lower Frame Rate videos.

drones, and robots are increasingly coming with higher frame rate cameras (240 FPS now being standard on many smart phones, tablets, drones, etc.). The visual object tracking community is yet to adapt to the changing landscape of what "real-time" means and how faster frame rates effect the choice of tracking algorithm one should employ.

In recent years, significant attention has been paid to Correlation Filter (CF) based methods [3, 13, 8, 15, 16, 1] for visual tracking. The appeal of correlation filters is their efficiency - a discriminative tracker can be learned online from a single frame and adapted after each subsequent

---

*Kiani & Fagg are joint first authors.
Email: hamedkg@gmail.com & ashton@fagg.id.au

frame. This online adaptation process allows for impressive tracking performance from a relatively low capacity learner. Further, CFs take advantage of intrinsic computational efficiencies afforded from operating in the Fourier domain [20]. Some CF methods (such as [13, 3]) are able to operate at hundreds of frames per second on embedded devices.

More recently, however, the visual tracking community has started to focus upon improving reliability and robustness through advances in deep learning [2, 31, 33]. While such methods have been shown to work well, their use comes at a cost. First, extracting discriminative features from CNNs or applying deep tracking frameworks is computationally expensive. Some deep methods operate at only a fraction of a frame per second, or require a high-end GPU to achieve real time performance. Second, training deep trackers can sometimes require a very large amount of data, as the learners are high capacity and require a significant amount of expense to train.

It is well understood that the central artifact that effects a visual tracker's performance is tolerance to appearance variation from frame to frame. Deep tracking methods have shown a remarkable aptitude for providing such tolerance - with the unfortunate drawback of having a considerable computational footprint [31, 33]. In this paper we want to explore an alternate thesis. Specifically, we want to explore that if we actually increase the frame rate - thus reducing the amount of appearance variation per frame - could we get away with substantially simpler tracking algorithms (from a computational perspective)? Further, could these computational savings be dramatic enough to cover the obvious additional cost of having to process a significant number of more image frames per second. Due to the widespread availability of higher frame rate (e.g. 240 FPS) cameras on many consumer devices we believe the time is ripe to explore such a question.

Inspired by the recent work of Handa et al. [11] in visual odometry, we believe that increasing capture frame rate allows for a significant boost in tracking performance without the need for deep features or complex trackers. We do not dismiss deep trackers or deep features, however we show that under some circumstances they are not necessary. By trading tracker capacity for higher frame rates - as is possible in many consumer devices - we believe that more favourable runtime performance can be obtained, particularly on devices with resource constraints, while still obtaining competitive tracking accuracy.

**Contributions:** In this paper, we present the Need for Speed (NfS) dataset and benchmark, the first benchmark (to our knowledge) for higher frame rate general object tracking using consumer devices. We use our dataset to evaluate numerous state of the art tracking methods (both CFs and deep learning based methods). An exciting outcome of our work was the unexpected result that if a sufficiently higher frame rate can be attained CFs with cheap hand-crafted features (*e.g.* HOG [5]) can outperform state of the art deep trackers in terms of accuracy and computational efficiency.

## 2. Related Work

### 2.1. Tracking Datasets

Standard tracking datasets such as OTB100 [37], VOT14 [17] and ALOV300 [30] have been widely used to evaluate current tracking methods in the literature. These datasets display annotated generic objects in real-world scenarios captured by low frame rate cameras (*i.e.* 24-30 FPS). Existing tracking datasets are briefly described as below.

**OTB50 and OTB100:** OTB50 [36] and OTB100 [37] belong to the Object Tracking Benchmark (OTB) with 50 and 100 sequences, respectively. OTB50 is a subset of OTB100, and both datasets are annotated with bounding boxes as well as 11 different attributes such as illumination variation, scale variation and occlusion and deformation.

**Temple-Color 128 (TC128):** This dataset consists of 128 videos which was specifically designed for the evaluation of color-enhanced trackers. Similar to OTBs, TC128 provides per frame bounding boxes and 11 per video attributes [22].

**VOT14 and VOT15:** VOT14 [17] and VOT15 [18] consist of 25 and 30 challenging videos, respectively, which are mainly borrowed from OTB100. All videos are labelled with rotated bounding boxes rather than upright ones. Both datasets come with per frame attribute annotation.

**ALOV300:** This dataset [29] contains 314 sequences mainly borrowed from the OTBs, VOT challenges and TC128. Videos are labeled with 14 visual attributes such as low contrast, long duration, confusion, zooming camera, motion smoothness, moving camera and transparency.

**UAV123:** This dataset [25] is recently created for Unmanned Aerial Vehicle (UAV) target tracking. There are 128 videos in this dataset, 115 videos captured by UAV cameras and 8 sequences rendered by a UAV simulator, which are all annotated with bounding boxes and 12 attributes.

Table 1 compares the NfS dataset with these datasets, showing that NfS is the only dataset with higher frame rate videos captured at 240 FPS. Moreover, in terms of number of frames, NfS is the largest dataset with 380K frames which is more than two times bigger than ALOV300.

### 2.2. Tracking Methods

Recent trackers can be generally divided into two categories, including correlation filter (CF) trackers [1, 13, 7, 23, 9] and deep trackers [26, 2, 34, 31]. We briefly review each of these two categories as following.

Table 1. Comparing NfS with other object tracking datasets.

| | UAV123 [25] | OTB50 [36] | OTB100 [37] | TC128 [22] | VOT14 [17] | VOT15 [18] | ALOV300 [29] | NfS |
|---|---|---|---|---|---|---|---|---|
| Capture frame rate | 30 | 30 | 30 | 30 | 30 | 30 | 30 | 240 |
| # Videos | 123 | 51 | 100 | 129 | 25 | 60 | 314 | 100 |
| Min frames | 109 | 71 | 71 | 71 | 171 | 48 | 19 | 169 |
| Mean frames | 915 | 578 | 590 | 429 | 416 | 365 | 483 | 3830 |
| Max frames | 3085 | 3872 | 3872 | 3872 | 1217 | 1507 | 5975 | 20665 |
| Total frames | 112578 | 29491 | 59040 | 55346 | 10389 | 21871 | 151657 | 383000 |

**Correlation Filter Trackers:** The interest in employing CFs for visual tracking was ignited by the seminal MOSSE filter [3] with an impressive speed of ∼700 FPS, and the capability of online adaptation. Thereafter, several works [13, 1, 6, 10, 24] were built upon the MOSSE showing notable improvement by learning CF trackers from more discriminative multi-channel features (*e.g.* HOG [5]) rather than pixel values. KCF [13] significantly improved MOSSE's accuracy by real-time learning of kernelized CF trackers on HOG features. Trackers such as Staple [1], LCT [24] and SAMF [10] were developed to improve KCF's robustness to object deformation and scale change. Kiani *et al.* [16] showed that learning such trackers in the frequency domain is highly affected by boundary effects, leading to suboptimal performance [8]. The CF with Limited Boundaries (CFLB) [16], Spatially Regularized CF (SRDCF) [8] and the Background-Aware CF (BACF) [14] have proposed solutions to mitigate these boundary effects in the Fourier domain, with impressive results.

Recently, learning CF trackers from deep Convolutional Neural Networks (CNNs) features [28, 19] has offered superior results on several standard tracking datasets [9, 7, 23]. The central tenet of these approaches (such as HCF [23] and HDT [27]) is that- even by per frame online adaptation- hand-crafted features such as HOG are not discriminative enough to capture the visual difference between consecutive frames in low frame rate videos. Despite their notable improvement, the major drawback of such CF trackers is their intractable complexity (∼0.2 FPS on CPUs) mainly needed for extracting deep features and computing Fourier transforms on hundreds of feature channels.

**Deep Trackers:** Recent deep learning based trackers [35, 26, 33] represent a new paradigm in tracking. Instead of hand-crafted features, a deep network trained for a non-tracking task (such as object recognition [19]) is updated with video data for generic object tracking. Unlike the direct combination of deep features with the traditional shallow methods *e.g.* CFs [9, 23, 7, 27], the updated deep trackers aim to learn from scratch the target-specific features for each new video. For example, the MDNet [26] learns generic features on a large set of videos and updates the so called domain-specific layers for unseen ones. The more re-

lated training set and unified training and testing approaches make MDNet win the first place in the VOT15 Challenge. Wang *et al.* [33] proposed to use fully convolutional networks (FCNT) with feature map selection mechanism to improve performance. However, such methods are computationally very expensive (even with a high end GPU) due to the fine-tuning step required to adapt the network from a large number of example frames. There are two high-speed deep trackers GOTURN [12] and SFC [2] that are able to run at 100 FPS and 75 FPS respectively on GPUs. Both of these methods train a Siamese network offline to predict motion between two frames (either using deep regression or a similarity comparison). At test time, the network is evaluated without any fine-tuning. Thus, these trackers are significantly less expensive because the only computational cost is the fixed feed-forward process. For these trackers, however, we remark that there are two major drawbacks. First, their simplicity and fixed-model nature can lead to high speed, but also lose the ability to update the appearance model online which is often critical to account for drastic appearance changes. Second, on modern CPUs, their speed becomes no more than 3 FPS, which is too slow for practical use on devices with limited computational resources.

## 3. NfS Dataset

The NfS datset consists of 100 higher frame rate videos captured at 240 FPS. We captured 75 videos using the iPhone 6 (and above) and the iPad Pro which are capable of capturing 240 frames per second. We also included 25 sequences from YouTube, which were captured at 240 FPS from a variety of different devices. All 75 captured videos come with corresponding IMU and Gyroscope raw data gathered during the video capturing process. Although we make no use of such data in this paper, we will make the data publicly available for potential applications.

The tracking targets include (but not limited to) vehicle (bicycle, motorcycle, car), person, face, animal (fish, bird, mammal, insect), aircraft (airplane, helicopter, drone), boat, and generic objects (*e.g.* sport ball, cup, bag, etc.). Each frame in NfS is annotated with an axis aligned bounding box using the VATIC toolbox [32]. Moreover, all videos are labeled with nine visual attributes, including occlusion, illu-

Table 2. Distribution of visual attributes within the NfS dataset, showing the number of coincident attributes across all videos. Please refer to Section 3 for more details.

|      | IV | SV | OCC | DEF | FM | VC | OV | BC | LR |
|------|----|----|-----|-----|----|----|----|----|----|
| IV   | **45** | 39 | 23 | 12 | 33 | 17 | 9  | 16 | 8  |
| SV   | 39 | **83** | 41 | 36 | 57 | 41 | 21 | 28 | 8  |
| OCC  | 23 | 41 | **51** | 21 | 31 | 23 | 10 | 19 | 8  |
| DEF  | 12 | 36 | 21 | **38** | 19 | 30 | 4  | 13 | 1  |
| FM   | 33 | 57 | 31 | 19 | **70** | 31 | 20 | 24 | 4  |
| VC   | 17 | 41 | 23 | 30 | 31 | **46** | 9  | 16 | 0  |
| OV   | 9  | 21 | 10 | 4  | 20 | 9  | **22** | 10 | 4  |
| BC   | 16 | 28 | 19 | 13 | 24 | 16 | 10 | **36** | 6  |
| LR   | 8  | 8  | 8  | 1  | 4  | 0  | 4  | 6  | **10** |

mination variation (IV), scale variation (SV), object deformation (DEF), fast motion (FM), viewpoint change (VC), out of view (OV), background clutter (BC) and low resolution (LR). The distribution of these attributes for NfS is presented in Table 2. Example frames of the NfS dataset and detailed description of each attribute are provided in the supplementary material.

## 4. Evaluation

**Evaluated Algorithms:** We evaluated 15 recent trackers on the NfS dataset. We generally categorised these trackers based on their learning strategy and utilized feature in three classes including CF trackers with hand-crafted features (BACF [14], SRDCF [8], Staple [1], DSST [6], KCF [13], LCT [24], SAMF [21] and CFLB [16]), CF trackers with deep features (HCF [23] and HDT [27]) and deep trackers (MDNet [26], SiameseFc [2], FCNT [33], GOTURN [12]). We also included MEEM [38] in the evaluation as the state of the art SVM-based tracker with hand-crafted feature. All these trackers are detailed in the supplementary material in terms of learning strategy and feature representation.

**Evaluation Methodology:** We use the success metric to evaluate all the trackers [36]. Success measures the intersection over union (IoU) of predicted and ground truth bounding boxes. The success plot shows the percentage of bounding boxes whose IoU is larger than a given threshold. We use the Area Under the Curve (AUC) of success plots to rank the trackers. We also compare all the trackers by their success rate at the conventional thresholds of 0.50 (IoU > 0.50) [36]. Moreover, we report the relative accuracy improvement which is computed as $\frac{\text{improved accuracy}}{\text{accuracy of lower frame rate tracking}}$, where the improved accuracy is the difference between accuracy (success rate at IoU > 0.50) of higher frame rate and lower frame rate tracking.

**Tracking Scenarios:** To measure the effect of capture frame rate on tracking performance, we consider two different tracking scenarios. At the first scenario, we run each
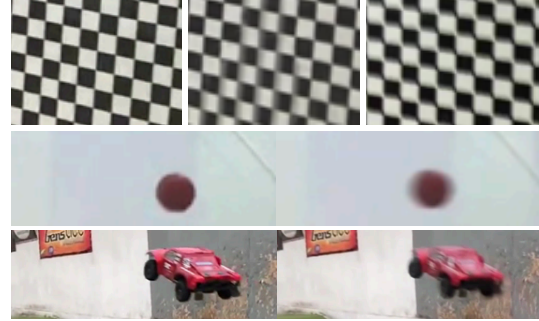


Figure 2. Top) a frame captured by a high frame rate camera (240 FPS), the same frame with synthesized motion blur, and the same frame captured by a low frame rate camera (30 FPS) with real motion blur. Bottom) sampled frames with corresponding synthesized motion blur. Please refer to Tracking Scenarios for more details.

tracker over all frames of the higher frame rate videos (240 FPS) in the NfS dataset. The second scenario, on the other hand, involves tracking lower frame rate videos (30 FPS). Since all videos in the NfS dataset are captured by high frame rate cameras, and thus no 30 FPS video is available, we simply create a lower frame rate version of NfS by temporal sampling every 8th frame. In such case, we track the object over each 8th frame instead of all frames. This simply models the large visual difference between two following sampled frames as one may observe in a real lower frame rate video. However, the main issue of temporal sampling is that since the videos are originally captured by higher frame rate cameras with very short exposure time, the motion blur caused by fast moving object/camera is significantly diminished. This leads to excluding the effect of motion blur in lower frame rate tracking. To address this concern and make the evaluation as realistic as possible, we simulate motion blur over the lower frame rate videos created by temporal sampling. We utilize a leading visual effects package (Adobe After Effects) to synthesize motion blur over the sampled frames. To verify the realism of the synthesized motion blur, Fig. 2 demonstrates a real frame (of a checkerboard) captured by a 240 FPS camera, the same frame with synthesized motion blur and a frame with real motion blur captured by a 30 FPS camera with identical extrinsic and intrinsic settings. To capture two sequences with different frame rates, we put two iPhones capture rates of 30 and 240 FPS side-by-side and then capture sequences from the same scene simultaneously. Fig. 2 also shows two frames before and after adding synthesized motion blur.

### 4.1. Adjusting Learning Rate of CF Trackers

A unique characteristic of CF trackers is their inherent ability to update the tracking model online, when new frames become available. The impact of each frame in learning/updating process is controlled by a constant weight

Table 3. Evaluating the effect of updating learning rate of each CF tracker on tracking higher frame rate videos (240 FPS). Accuracy is reported as success rate (%) at IoU > 0.50. Please refer to Section 4.1 for more details about the original and updated learning rates.

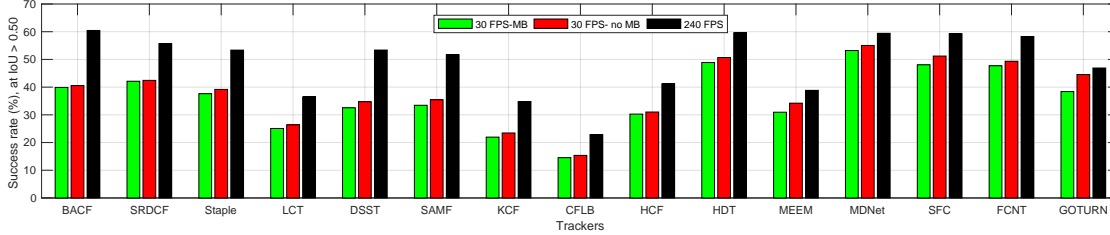|  | BACF | SRDCF | Staple | LCT | DSST | SAMF | KCF | CFLB | HCF | HDT |
|---|---|---|---|---|---|---|---|---|---|---|
| Original LR | 48.8 | 48.2 | 51.1 | 34.5 | 44.0 | 42.8 | 28.7 | 18.3 | 33.0 | 57.7 |
| Updated LR | 60.5 | 55.8 | 53.4 | 36.4 | 53.4 | 51.7 | 34.8 | 22.9 | 41.2 | 59.6 |



Figure 3. Comparing higher frame rate tracking (240 FPS) versus lower frame rate tracking (30 FPS) for each tracker. For higher frame rate tracking CF trackers employ updated learning rates. The results of lower frame rate tracking are plotted for videos with and without motion blur (30 FPS-MB and 30 FPS- no MB). Results are reported as success rate (%) at IoU > 0.50.

called learning (or adaptation) rate [3]. Smaller rates increase the impact of older samples, while bigger ones give higher weight to samples from more recent frames. Each CF tracker has its own learning rate which was tuned for robust tracking on low frame rate sequences (30 FPS). To retain the robustness of such methods for higher frame rate videos (240 FPS), we approximately adjust their learning rate to be $LR_{new} = \frac{1}{8}LR_{old}$. Since the number of frames in 240 FPS videos is 8 times more than that in 30 FPS sequences over a fixed period of time, dividing learning rates by 8 can keep the balance between CFs updating capacity and smaller inter-frame variation in 240 FPS videos [1].

Here, we empirically demonstrate how adjusting the learning rates of CF trackers affects their tracking performance. Table 3 shows the tracking accuracy (success rate at IoU > 0.50) of 10 recent CF trackers on the NfS 240 FPS sequences, comparing tracking by original learning rates ($LR_{old}$ from their reference papers) versus updated rates $LR_{new}$. The result shows that adjusting the learning rates notably improves the accuracy of all the CF trackers. For some trackers such as BACF, SRDCF, DSST and HCF there is a substantial improvement, while for Staple, LCT and HDT the improvement is much smaller ($\sim 2\%$). This is most likely because of the complementary parts of these trackers. Staple utilizes color scores per pixel, and LCT uses random fens classifiers as additional detectors/trackers which are independent of their CF modules. Similarly, HDT employs the Hedge algorithm [4] as a multi-experts decision maker to merge hundreds of weak CF trackers in a strong tracker. Thus, updating their learning rates offers less improvement compared to those trackers such as BACF and DSST that solely track by a single CF based tracker.

---

[1]Proof is provided in the supplementary material.

## 4.2. Per Tracker Evaluation

Figure 3 compares tracking higher versus lower frame rate videos for each evaluated method. For lower frame rate tracking (30 FPS) results are reported for both with and without motion blur. All CF trackers achieve a significant increase in performance (AUCs are improved > 10%) when tracking on 240 FPS videos. This is because in higher frame rate video, the appearance change between two adjacent frames is very small, which can be effectively learned by per frame CF online adaptation. Among deep trackers, FCNT achieved the most improvement (6%), since this tracker also fine-tunes every 20 frames using the most confident tracking results. The lowest improvement belongs to SFC and MD-Net. These methods are trained off-line and do not update a model or maintain a memory of past appearances [2]. Thus, tracking higher frame rate videos offers much smaller improvement to such trackers. When evaluating lower frame rate videos, a slight performance drop can be observed with the presence of motion blur, demonstrating that all trackers are reasonably robust towards motion blur.

## 4.3. Overall Comparison

The overall comparison of all trackers over three tracking settings- higher frame rate tracking (240 FPS), lower frame rate tracking with synthesized motion blur (30 FPS MB) and lower frame rate tracking without motion blur (30 FPS no MB)- is demonstrated in Fig. 4 (success plots) and Table 4 (AUCs and tracking speed).

**Accuracy Comparison:** For lower frame rate tracking without motion blur MDNet achieved the best performance followed by SFC. HDT which utilizes deep features over a CF framework obtained the third rank followed by FCNT. Almost the same ranking is observed for lower frame rate tracking with motion blur. Overall, deep trackers outper-
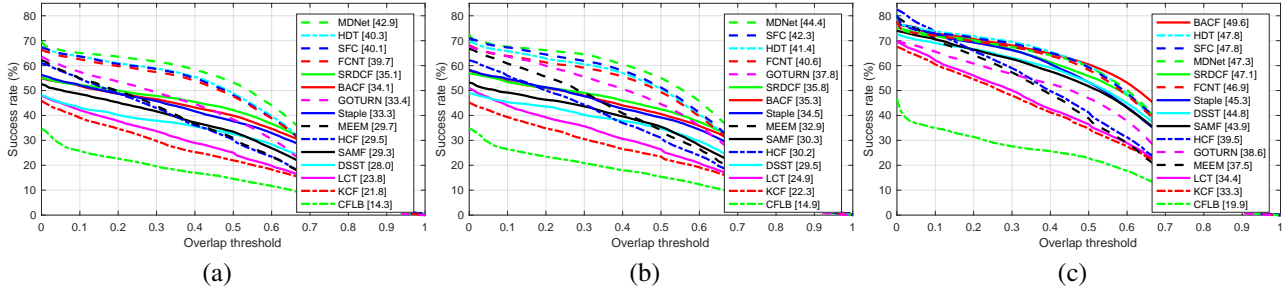
Figure 4. Evaluating trackers over three tracking scenarios, (a) lower frame rate tracking with synthesized motion blur, (b) and lower frame rate tracking without motion blur, and (c) higher frame rate tracking. AUCs are reported in brackets.

Table 4. Comparing trackers on three tracking scenarios including higher frame rate tracking (240 FPS), lower frame rate tracking with synthesized motion blur (30 FPS MB) and lower frame rate tracking without motion blur (30 FPS no MB). Results are reported as the AUC of success plots. We also show the speed of each tracker on CPUs and/or GPUs if applicable. The first, second, third, forth and fifth highest AUCs/speeds are highlighted in color.

| | BACF | SRDCF | Staple | LCT | DSST | SAMF | KCF | CFLB | HCF | HDT | MEEM | MDNet | SFC | FCNT | GOTURN |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 30 FPS- no MB | 35.2 | 35.7 | 34.5 | 24.8 | 29.4 | 30.3 | 22.3 | 14.9 | 30.2 | 41.3 | 32.9 | 44.4 | 42.3 | 40.5 | 37.7 |
| 30 FPS- MB | 34.0 | 35.1 | 33.2 | 23.7 | 28.0 | 29.2 | 21.7 | 14.2 | 29.5 | 40.3 | 29.6 | 42.9 | 40.1 | 39.7 | 33.4 |
| 240 FPS | 49.5 | 47.1 | 45.3 | 34.3 | 44.8 | 43.9 | 33.3 | 19.9 | 39.5 | 47.8 | 37.5 | 47.3 | 47.7 | 46.9 | 38.6 |
| Speed (CPU) | 38.3 | 3.8 | 50.8 | 10.0 | 12.5 | 16.6 | 170.4 | 85.1 | 10.8 | 9.7 | 11.1 | 0.7 | 2.5 | 3.2 | 3.9 |
| Speed (GPU) | - | - | - | - | - | - | - | - | - | 43.1 | - | 2.6 | 48.2 | 51.8 | 155.3 |

formed CF trackers for lower frame rate tracking. This is not surprising, as deep trackers have a high learning capacity and employ highly discriminative deep features which are able to handle the large variation in adjacent frames which is present in lower frame rate videos.

Surprisingly, the best accuracy of higher frame rate tracking achieved by BACF (49.56), which is a CF tracker with HOG features, followed by HDT (47.80)- a CF tracker with deep features. SRDCF (47.13), Staple (45.34), DSST (44.80) and SAMF (43.92) outperformed GOTURN (38.65) and obtained very competitive accuracy compared to other deep trackers including SFC (47.78), MDNet (47.34) and FCNT (46.94). This implies that when higher frame rate videos are vailable, the ability of CF trackers to adapt online is of greater benefit than high learning capacity of deep trackers. The reasoning for this is intuitive, since for higher frame rate video there is less appearance change among consecutive frames, which can be efficiently modeled by updating the tracking model at each frame even using simple hand-crafted features.

**Run-time Comparison:** The tracking speed of all evaluated methods in FPS is reported in Table 4. For the sake of fair comparison, we tested MATLAB implementations of all methods (including deep trackers) on a 2.7 GHz Intel Core i7 CPU with 16 GB RAM. We also reported the speed of deep trackers on nVidia GeForce GTX Titan X GPU to have a better sense of their run-time when GPUs are available. On CPUs, all CF trackers achieved much higher speed compared to all deep trackers, because of their shallow ar-

chitecture and efficient computation in the Fourier domain. Deep trackers on CPUs performed much slower than CFs, with the exception of SRDCF (3.8 FPS). On GPU, however, deep trackers including GOTURN (155.3 FPS), FCNT (51.8 FPS) and SFC (48.2 FPS) performed much faster than on CPU. Their performance is comparable with many CF trackers running on CPU, such as KCF (170.4 FPS), BACF (38.3 FPS) and Staple (50.8 FPS). For tracking lower frame rate videos, only BACF, Staple, KCF and CFLB can track at or above real time on CPUs. GOTURN, FCNT and SFC offer real-time tracking of lower frame rate videos on GPUs. KCF (170.4 FPS) and GOTURN (155.3 FPS) are the only trackers which can track higher frame rate videos almost real time on CPUs and GPUs, respectively.

### 4.4. Attribute-based Evaluation

The attribute based evaluation of all trackers on three tracking settings is shown in Fig. 6 (success rate at IoU > 0.50) and Table 5 (relative accuracy improvement). Similar to the previous evaluation, the presence of motion blur slightly degrades the performance of tracking lower frame rate sequences. For lower frame rate tracking, in general, deep trackers outperformed CF trackers over all nine attributes, MDNet outperformed SFC and FCNT, and HDT achieved the superior performance compared to other CF trackers for all attributes. Compared to lower frame rate tracking, tracking higher frame rate videos offers a notable improvement of all trackers on all attributes except non-rigid deformation, as can be seen in Table 5. This demon-
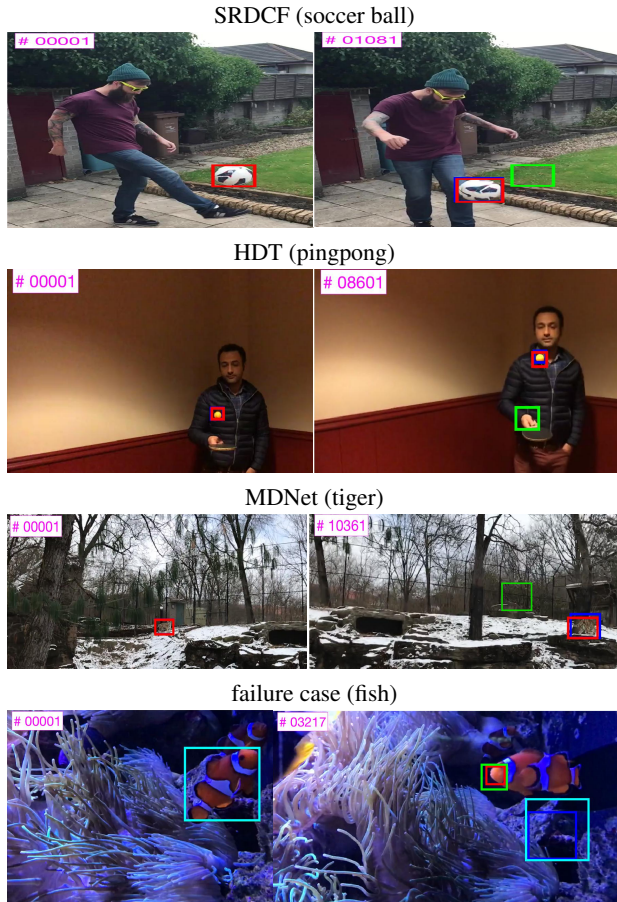
SRDCF (soccer ball)

HDT (pingpong)

MDNet (tiger)

failure case (fish)

Figure 5. Rows(1-3) show tracking performance of three trackers icluding a CF tracker with HOG (SRDCF), a CF tracker with deep features (HDT) and a deep tracker (MDNets), comparing lower frame rate (green boxes) versus higher frame rate (red boxes) tracking. Ground truth is shown by blue boxes. Last row visualizes a failure case of higher frame rate tracking caused by non-rigid deformation for BACF, Staple, MDNet and SFC.

strates the sensitivity of both CF based and deep trackers to non-rigid deformation even when they track higher frame rate videos. Fig. 6 shows that CF trackers with hand-crafted features outperformed all deep trackers as well as HDT for 6 attributes. More particularly, for illumination variation, occlusion, fast motion, out-of-view, background clutter and low resolution CF trackers with hand-crafted features (such as BACF and SRDCF) achieved superior performance to all deep trackers and HDT. However, deep tracker MDNet achieved the highest accuracy for scale variation (61.0), deformation (59.2) and view change (55.9), closely followed by BACF for scale variation (60.1) and HDT for deformation (57.6) and view change (54.8). The relative accuracy improvement of tracking higher frame rate versus lower frame rate videos (with motion blur) for each tracker and each attribute is reported in Table 5. The result shows that

Table 5. Relative accuracy improvement (%) of high frame rate tracking versus low frame rate tracking for each attribute. Relative improvement more than 50% are underlined in red.

|  | IV | SV | OCC | DEF | FM | VC | OV | BC | LR |
|---|---|---|---|---|---|---|---|---|---|
| BACF | 82.1 | 44.5 | 50.8 | 11.7 | 62.1 | 26.1 | 39.6 | 47.2 | 118.3 |
| SRDCF | 28.4 | 19.1 | 40.7 | 11.3 | 43.9 | 36.9 | 13.0 | 59.3 | 47.2 |
| Staple | 106.0 | 37.1 | 42.7 | 10.3 | 55.9 | 22.5 | 33.8 | 11.6 | 66.9 |
| LCT | 74.5 | 22.4 | 60.7 | 16.3 | 48.7 | 36.2 | 18.3 | 94.8 | 127.8 |
| DSST | 119.4 | 58.4 | 57.9 | 12.1 | 84.5 | 41.1 | 36.0 | 65.3 | 54.4 |
| SAMF | 124.4 | 49.9 | 62.6 | 13.2 | 73.0 | 28.2 | 33.1 | 52.6 | 44.9 |
| KCF | 128.0 | 37.7 | 66.7 | 5.1 | 78.7 | 34.3 | 21.4 | 81.9 | 123.6 |
| CFLB | 113.2 | 66.7 | 52.3 | 5.4 | 79.1 | 14.0 | 42.3 | 89.3 | 200.4 |
| HCF | 50.1 | 23.1 | 62.8 | 16.2 | 39.0 | 26.0 | 30.6 | 71.8 | 134.5 |
| HDT | 29.2 | 19.1 | 12.5 | 4.3 | 26.2 | 16.6 | 19.8 | 13.8 | 32.5 |
| MEEM | 32.2 | 17.8 | 33.6 | 16.3 | 31.1 | 19.1 | 14.9 | 29.0 | 60.9 |
| MDNet | 14.1 | 8.3 | 9.7 | 1.9 | 16.5 | 13.9 | 4.1 | 16.9 | 10.6 |
| SFC | 35.4 | 20.1 | 20.6 | 7.9 | 26.7 | 16.6 | 15.7 | 22.1 | 27.8 |
| FCNT | 29.7 | 19.3 | 16.1 | 4.8 | 24.8 | 19.1 | 15.5 | 12.5 | 26.5 |
| GOTURN | 52.9 | 27.7 | 16.5 | -3.8 | 38.3 | 18.2 | 34.9 | 37.2 | -8.6 |

that first, compared to other attributes, less improvement achieved for non-rigid deformation attribute for all trackers, and second, the percentage of relative improvement for CF trackers is much higher than that of deep trackers.

### 4.5. Discussion and Conclusion

In this paper, we introduce the first higher frame rate object tracking dataset and benchmark. We empirically evaluate the performance of the state-of-the-art trackers with respect to two different capture frame rates (30 FPS vs. 240 FPS), and find the surprising result that at higher frame rates, simple trackers such as correlation filters trained on hand-crafted features (e.g. HOG) outperform complex trackers based on deep architecture. This suggests that computationally tractable methods such as cheap CF trackers in conjunction with higher capture frame rate videos can be utilized to effectively perform object tracking on devices with limited processing resources such as smart phones, tablets, drones, etc. As shown in Fig. 7, cheaper trackers on higher frame rate video (e.g. KCF and Staple) have been demonstrated as competitive with many deep trackers on lower frame rate videos (such as HDT and FCNT).

Our results also suggest that traditional evaluation criteria that trades off accuracy versus speed (e.g., Fig.7 in [18]) could possibly paint an incomplete picture. This is because, up until now, accuracy has been measured without regard to the frame rate of the video. As we show, this dramatically underestimates the performance of high speed algorithms. In simple terms: the accuracy of a 240 FPS tracker cannot be truly appreciated until it is run on a 240 FPS video! From an embedded-vision perspective, we argue that the acquisition frame rate is a resource that should be explicitly traded off when designing systems, just as is hardware (GPU vs CPU). Our new dataset allows for, the first time, exploration
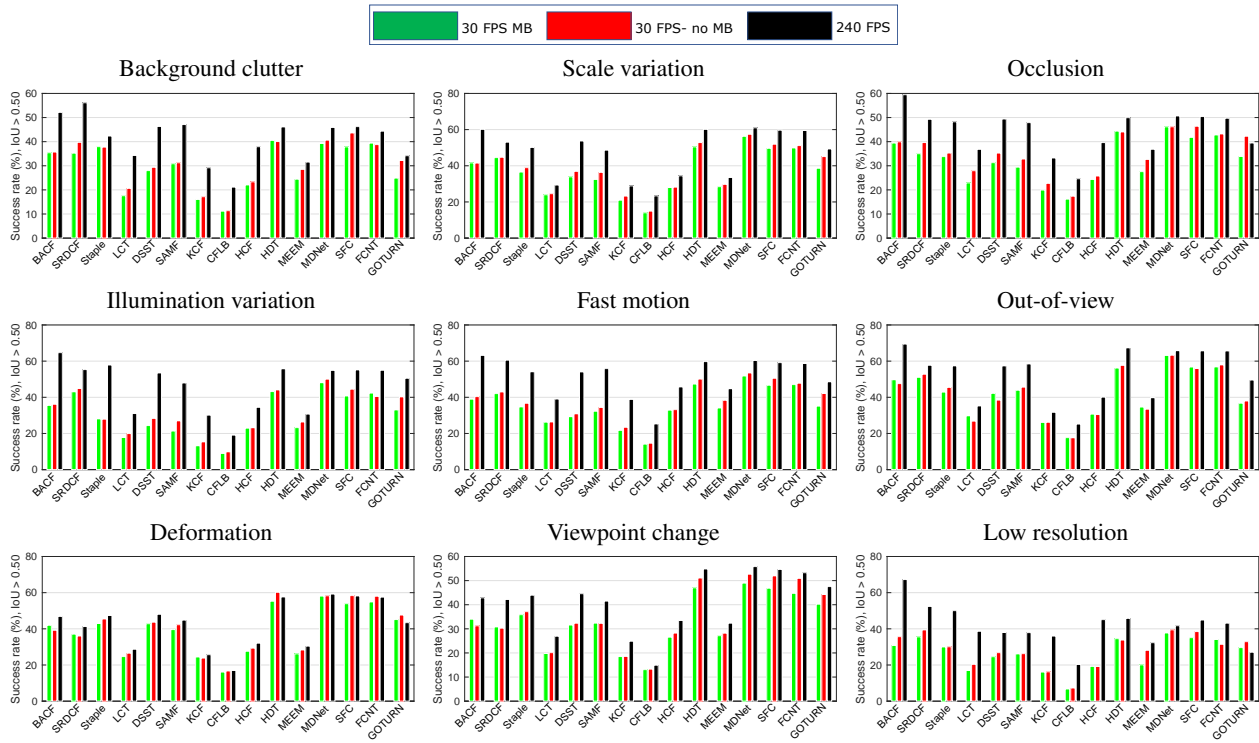
Figure 6. Attribute based evaluation. Results are reported as success rate (%) at IoU > 0.50.
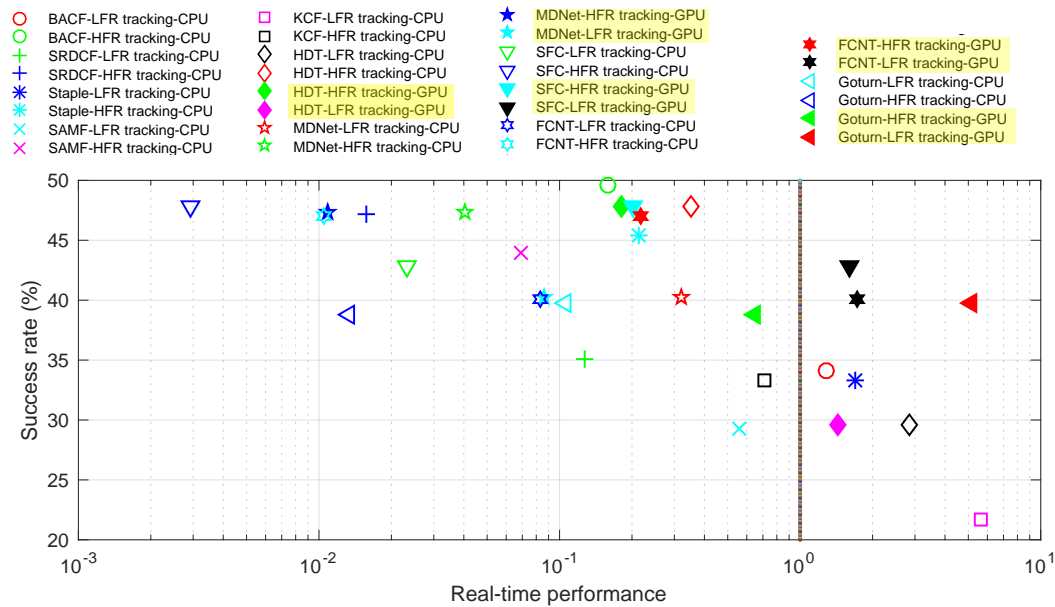


Figure 7. This plot shows the affect of resource availability (GPUs vs. CPUs) and frame rate of captured videos (lower vs. higher frame rate) on the top-10 evaluated trackers' accuracies and real-time performance. Real-time performance is computed as the ratio of each tracker's speed (FPS) to frame rate of the target videos (30 vs. 240 FPS). The vertical line on the plot shows the boundary of being real-time (frame rate of the target video is the same as the tracker's speed). Trackers which are plotted at the left side of the line are not able to track real-time (according to their tracking speed and video frame rate). GPU results are highlighted in yellow.

of such novel perspectives. Our dataset fills a need. The need for speed.

# References

[1] L. Bertinetto, J. Valmadre, S. Golodetz, O. Miksik, and P. H. S. Torr. Staple: Complementary learners for real-time tracking. In *CVPR*, June 2016. 1, 2, 3, 4

[2] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. Torr. Fully-convolutional siamese networks for object tracking. In *ECCV*, pages 850–865, 2016. 1, 2, 3, 4, 5

[3] D. S. Bolme, J. R. Beveridge, B. A. Draper, and Y. M. Lui. Visual object tracking using adaptive correlation filters. In *CVPR*, pages 2544–2550. IEEE, 2010. 1, 2, 3, 4

[4] K. Chaudhuri, Y. Freund, and D. J. Hsu. A parameter-free hedging algorithm. In *NIPS*, pages 297–305, 2009. 5

[5] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *CVPR*, volume 1, pages 886–893, 2005. 2, 3

[6] M. Danelljan, G. Häger, F. Khan, and M. Felsberg. Accurate scale estimation for robust visual tracking. In *BMVC*, 2014. 3, 4

[7] M. Danelljan, G. Hager, F. Shahbaz Khan, and M. Felsberg. Convolutional features for correlation filter based visual tracking. In *CVPR Workshop*, pages 58–66, 2015. 2, 3

[8] M. Danelljan, G. Hager, F. Shahbaz Khan, and M. Felsberg. Learning spatially regularized correlation filters for visual tracking. In *ICCV*, pages 4310–4318, 2015. 1, 3, 4

[9] M. Danelljan, A. Robinson, F. S. Khan, and M. Felsberg. Beyond correlation filters: Learning continuous convolution operators for visual tracking. In *ECCV*, pages 472–488, 2016. 2, 3

[10] M. Danelljan, F. Shahbaz Khan, M. Felsberg, and J. Van de Weijer. Adaptive color attributes for real-time visual tracking. In *CVPR*, pages 1090–1097, 2014. 1, 3

[11] A. Handa, R. A. Newcombe, A. Angeli, and A. J. Davison. Real-time camera tracking: When is high frame-rate best? In *ECCV*, pages 222–235, 2012. 2

[12] D. Held, S. Thrun, and S. Savarese. Learning to track at 100 fps with deep regression networks. In *ECCV*, 2016. 3, 4

[13] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista. High-speed tracking with kernelized correlation filters. *PAMI*, 37(3):583–596, 2015. 1, 2, 3, 4

[14] H. Kiani Galoogahi, A. Fagg, and S. Lucey. Learning background-aware correlation filters for visual tracking. *arXiv preprint 1703.04590*, 2017. 3, 4

[15] H. Kiani Galoogahi, T. Sim, and S. Lucey. Multi-channel correlation filters. In *ICCV*, pages 3072–3079, 2013. 1

[16] H. Kiani Galoogahi, T. Sim, and S. Lucey. Correlation filters with limited boundaries. In *CVPR*, pages 4630–4638, 2015. 1, 3, 4

[17] M. Kristan, J. Matas, A. Leonardis, and et al. The visual object tracking vot2014 challenge results. In *ECCV Workshop*, 2014. 2, 3

[18] M. Kristan, J. Matas, A. Leonardis, M. et alFelsberg, L. Cehovin, G. Fernandez, T. Vojir, G. Hager, G. Nebehay, and R. Pflugfelder. The visual object tracking vot2015 challenge results. In *CVPR Workshop*, pages 1–23, 2015. 2, 3, 7

[19] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, pages 1097–1105, 2012. 3

[20] B. V. Kumar. *Correlation pattern recognition*. Cambridge University Press, 2005. 2

[21] Y. Li and J. Zhu. A scale adaptive kernel correlation filter tracker with feature integration. In *ECCV*, pages 254–265, 2014. 4

[22] P. Liang, E. Blasch, and H. Ling. Encoding color information for visual tracking: algorithms and benchmark. *ITIP*, 24(12):5630–5644, 2015. 1, 2, 3

[23] C. Ma, J.-B. Huang, X. Yang, and M.-H. Yang. Hierarchical convolutional features for visual tracking. In *CVPR*, pages 3074–3082, 2015. 2, 3, 4

[24] C. Ma, X. Yang, C. Zhang, and M.-H. Yang. Long-term correlation tracking. In *CVPR*, pages 5388–5396, 2015. 3, 4

[25] M. Mueller, N. Smith, and B. Ghanem. A benchmark and simulator for uav tracking. In *ECCV*, pages 445–461, 2016. 2, 3

[26] H. Nam and B. Han. Learning multi-domain convolutional neural networks for visual tracking. *arXiv preprint 1510.07945*, 2015. 2, 3, 4

[27] Y. Qi, S. Zhang, L. Qin, H. Yao, Q. Huang, J. Lim, and M.-H. Yang. Hedged deep tracking. In *CVPR*, pages 4303–4311, 2016. 3, 4

[28] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint 1409.1556*, 2014. 3

[29] A. W. Smeulders, D. M. Chu, R. Cucchiara, S. Calderara, A. Dehghan, and M. Shah. Visual tracking: An experimental survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7). 2, 3

[30] A. W. Smeulders, D. M. Chu, R. Cucchiara, S. Calderara, A. Dehghan, and M. Shah. Visual tracking: An experimental survey. *PAMI*, 36(7):1442–1468, 2014. 2

[31] R. Tao, E. Gavves, and A. W. Smeulders. Siamese instance search for tracking. *arXiv preprint 1605.05863*, 2016. 2

[32] C. Vondrick, D. Patterson, and D. Ramanan. Efficiently scaling up crowdsourced video annotation. *IJCV*, 101(1):184–204, 2013. 3

[33] L. Wang, W. Ouyang, X. Wang, and H. Lu. Visual tracking with fully convolutional networks. In *ICCV*, pages 3119–3127, 2015. 2, 3, 4

[34] L. Wang, W. Ouyang, X. Wang, and H. Lu. Stct: Sequentially training convolutional networks for visual tracking. CVPR, 2016. 2

[35] N. Wang and D.-Y. Yeung. Learning a deep compact image representation for visual tracking. In *NIPS*, 2013. 3

[36] Y. Wu, J. Lim, and M.-H. Yang. Online object tracking: A benchmark. In *CVPR*, pages 2411–2418, 2013. 1, 2, 3, 4

[37] Y. Wu, J. Lim, and M.-H. Yang. Object tracking benchmark. *PAMI*, 37(9):1834–1848, 2015. 1, 2, 3

[38] J. Zhang, S. Ma, and S. Sclaroff. MEEM: robust tracking via multiple experts using entropy minimization. In *ECCV*, pages 188–203, 2014. 4