Kourosh Meshgi, Shigeyuki Oba Graduate School of Informatics, Kyoto University

Abstract

Recently, discriminative visual trackers obtain state-of-the-art performance, yet they suffer in the presence of different real-world challenges such as target motion and appearance changes. In a discriminative tracker, one or more classifiers are employed to obtain the target/non-target label for the samples, which in turn determine the target's location. To cope with variations of the target shape and appearance, the classifier(s) are updated online with different samples of the target and the background. Sample selection, labeling and updating the classifier is prone to various sources of errors that drift the tracker. In this study we motivate, conceptualize, realize and formalize a novel active co-tracking framework, step-by-step to demonstrate the challenges and generic solutions for them. In this framework, not only classifiers cooperate in labeling the samples, but also exchange their information to robustify the labeling, improve the sampling, and realize efficient yet effective updating. The proposed framework is evaluated against state-of-the-art trackers on public dataset and showed promising results. keywords: visual tracking, active learning, active co-tracking, uncertainty sampling

²² 1 Introduction

2

3

8

q

10

11

12

13

14

15

16

17

18

19

20

21

Visual tracking is one of the building blocks of human-robot interaction. Implicit or
explicit, this task is embedded in many high-level complicated tasks of the robot. Attending the speaker in a multimodal spoken dialog system [1], following the target
[2], imitating the behavior of a human [3], extracting tacit information of an interaction [4], sign-language interpretation [5], and autonomous driving as well as simpler
tasks such as human-robot cooperation [6], obstacle avoidance [7], first-person view
action recognition [8] and human-computer interfaces [9].

The most general type of tracking is single-object model-free online tracking, in 30 which the object is annotated in the first frame, and tracked in the subsequent frames 31 with no prior knowledge about the target's appearance, its motions, the background, 32 the configurations of the camera, and other conditions of the scene. Visual tracking is 33 still considered as a challenging problem despite numerous efforts made to address 34 abrupt appearance changes of the target [10], its complex transformations [11] and 35 deformations [12], background clutter [13], occlusion [14], and motion artifacts [15]. 36 Generative trackers attempt to construct a robust object appearance model, or 37

to learn it on-the-fly using advanced machine learning techniques such as subspace

³⁹ learning [16], hash learning [17], dictionary learning [18], and sparse code learning

[10]. Although some settings allow for strong assumptions about the target, in real-40 world applications it is desired to track arbitrary objects with little a-priori knowl-41 edge. Such model free tracker consists of learning and adjusting the representation of 42 the target on-the-fly. To this end, discriminative models focus on target/background 43 separation using correlation filters [19, 20, 21] or dedicated classifiers [22], which 44 assist them to dominate the visual tracking benchmarks [23, 24, 25]. Using tracking-45 by-detection approaches is a popular trend in recent years, due to significant break-46 throughs in object detection domain (deep residual neural networks [26] for in-47 stance), yielding strong discriminating power with offline training. Adopted for vi-48 sual tracking, many of such trackers are adjusted for online training and accumulate 49 knowledge about a target with each successful detection (e.g., [27, 28, 29, 22]. 50

Tracking-by-detection methods primarily treat tracking as a detection problem to avoid having model object dynamics especially in the case of sudden motion changes, extreme deformations, and occlusions [30, 31]. However, there is a multitude of drawbacks in the tracking-by-detection setting:

Label noise: Inaccurate labels confuse the classifier [12] and degrade the classifier fication accuracy [30], The labeler is typically built upon heuristics and intuitions, rather than using the accumulated knowledge about the target.

58

59

2. *Self-learning loop*: the classifier is re-trained by their own output from earlier frames, thus accumulating error over time [31],

3. Uniform treatment of samples: Equal-weight for all samples in evaluating the target [32] and training the classifier [33], despite the uneven contextual information in different samples. The classifier is trained using all the examples with equal weights, meaning that negative examples which overlap very little with the target bounding box are treated equally as those negative examples with significant overlaps.

4. Stationarity assumption: Assuming a stationary distribution of the target appearance does not hold for most of the real-world scenarios with drastic target appearance changes [31]. In the context of visual tracking, the non-stationarity means that the appearance of an object may change so significantly that a negative sample in the current frame looks more similar to a positive example in the previous frames.

5. Model update difficulties: Adaptive trackers inherently suffer from the drifting
problem. Noisy model update [34] and the mismatch between model update
frequency and target evolution rate [35] are two major challenges of the model
update. if the update rate is small, the changes of the target are not reflected
into target's template, whereas rapid update of the tracker renders it vulnerable to data noise and small target localization errors. This phenomenon is also
known as stability plasticity dilemma.

In this study we motivate, conceptualize, realize and formalize a novel co-tracking framework. First, the importance of such system is demonstrated by a recent and comprehensive literature review. Then a discriminative tracking framework is formalized to be evolved to a co-tracking by explaining all the steps, mathematically

and intuitively. We then construct various instances of the proposed co-tracking

⁸⁴ framework (Table 1), to demonstrate how different topologies of the system can be

realized, how the information exchange is optimized, and how different challenges
 of tracking (e.g., abrupt motions, deformations, clutter) can be handled in the pro-

posed framework. Active learning will be explored in the context of labeling and

³⁸ information exchange of this co-tracking framework to speed up the tracker's con-

⁸⁹ vergence while updating the tracker's classifiers effectively. Dual memory is also

⁹⁰ proposed in the co-tracking framework to handle various tracking scenarios ranging from camera motions to temporal appearance changes of the target and occlusions.

Table 1: Trackers introduced in this chapter: **T0**: a part-based tracker without model update, **T1**: the part-based tracker with model update, **T2**: a KNN-based tracker with color and HOG features, **T3**: co-tracking of KNN-based classifier T2 and part-based detector T1, **T4**: active co-tracking of T1 and T2 with online update, **T5**: active asymmetric co-tracking of short-memory T1 and long-memory T2 (modified from [36]), and **T6**: active ensemble co-tracking of bagging-induced ensemble and long-memory T2 (modified from [37])

	T0	T1	T2	T3	T4	T5	T6
Online Update		\checkmark	\checkmark	~	~	~	\checkmark
Co-tracking				\checkmark	\checkmark	\checkmark	\checkmark
Active Learning					\checkmark	\checkmark	\checkmark
Dual Memory						\checkmark	\checkmark
Ensemble							\checkmark

91

⁹² 2 Tracking-by-detection

Typically tracking-by-detection method consists of five major steps: Sampling, Clas sifying, Labeling, Estimating, Updating.

95 SAMPLING: To obtain the positive sample(s) and negative samples (the target and

the background respectively), dense or sparse (stochastic) sampling is performed ei ther around last known target position (using Gaussian distributions, particle filters,

or various motion models) or around the saliencies or keypoints in the current frame

⁹⁹ [17]. Adaptive weights for the samples based on their appearance similarity to the

target [38], occlusion state [14], and spatial distance to previous target location [39]
 have been considered, especially in the context of tracking-by-detection, boosting
 [40] have been extensively investigated [41, 42, 43].

CLASSIFYING: The classification module of tracking-by-detection schemes utilizes offline-trained classifiers or online supervised learning methods to classify the target from its background (e.g. [44]). To robustify this module especially against label noise, supervised learning with robust loss functions [42, 45], semi-supervised

¹⁰⁷ [35, 46] and multi-instance [43, 47, 48] learning approaches are considered. Efficient

- ¹⁰⁸ sparse sampling [49], leveraging context information [50, 13], considering sample in-
- formation content for the classifier [51], and landmark-based label propagation [39]
 are among other proposed approaches to address this issue. Another interesting ap-

proach is to reformulate to couple the labeling and updating process to bridge the

gap between the objective of these two steps, as labeling aims for predicting binary sample labels whereas updating typically tries to estimate object location [12]. The label noise problem amplifies when the tracker does not have a forgetting mechanism or a way to obtain external scaffolds (i.e. self-learning loop). This inspired the use of co-tracking [30], ensemble tracking [52, 53] or label verification schemes [54] to break the self-learning loop using auxiliary classifiers.

LABELING: The result of classification process provides the target/background la bel for each sample, a process which can be enhanced by employing an ensemble of
 classifiers [52, 53], exchanging information between collaborative classifiers [30] and
 verifying labels by auxiliary classifiers [54] or landmarks [39].

ESTIMATING: The state of the target, i.e. the location and scale of the target usually
 described with a bounding box, is then determined by selecting the sample with the
 highest classification score [12], calculating the expectation of target state [37], or
 performing an estimated bounding box regression [55].

UPDATING: Updating the classifier is another challenge of the tracking-by-detection 126 schemes. Some researchers believe in the necessity of having a "teacher" to train the 127 classifier [35]. Adaptive ensemble of classifiers [56] and co-learning [30] in which 128 multiple trackers with different features or inference engines train each other aimed 129 to address this need using other detectors or trackers. Furthermore, some approaches 130 selected the most discriminative feature selection [41], combined generative and dis-131 criminative models [57], replaced the weakest classifier of an ensemble [41] or the 132 oldest one [56], or applied a budget on the sample pool of the classifier (hence, keep-133 ing only some prototypical samples) [12, 39], to overcome this problem. 134

On top of that, the frequency of update is another important role-player in tracker's 135 performance [35]. Higher update rates capture the rapid target changes, but is prone 136 to occlusions, whereas slower update paces provide a long memory for the tracker to 137 handle temporal target variations but lack the flexibility to accommodate permanent 138 target changes. To this end, researchers try to combine long- and short-term mem-130 ories [58], role-back improper updates [53], or utilize different temporal-snapshots 140 of the classifier to overcome non-stationary distribution of the target's appearance 141 [59]. This pipeline, however, was altered in some studies to introduce desired prop-142 erties, e.g., to avoid label noise by merging sampling and labeling steps [12]. 143

144 2.1 Formalization

Online visual tracking is the task to update the state vector \mathbf{p}_t involving location, size, and shape of the bounding box, at each observation of video frame t = 1, ..., T. The update process is sometimes written with transformation \mathbf{y}_t that transforms the previous state vector \mathbf{p}_{t-1} to the current state $\mathbf{p}_t = \mathbf{p}_{t-1} \circ \mathbf{y}_t$.

In tracking-by-discrimination framework, we utilize a classifier θ_t that discriminates an image patch **x** into either target or background, where the classifier is denoted as a real valued discriminant function $h(\mathbf{x}|\theta_t) \in \mathbb{R}$ and the function value $s = h(\mathbf{x}|\theta_t)$ is called a discrimination score, or in short, score. The patch **x** (i.e. the

- area of the image bounded by the bounding box \mathbf{p}_t) is labeled as target if $s > \tau$ with a threshold τ and as background if $x < \tau$. A typical procedure of the tracking-by-
- ¹⁵⁵ discrimination is written as follows.
- **SAMPLING**: The samples are defined using these transformations, and their corresponding image patches $\mathbf{x}_t^j \in \mathcal{X}_t$ are selected from image. We obtain N samples of state p_t^j , j = 1, ..., N by drawing random transformations $\mathbf{y}_t^j \in \mathcal{Y}_t$ using dense or sparse sampling strategy, transforming the previous state p_{t-1} with a transformations \mathbf{y}_t^j as $\mathbf{p}_t^j = \mathbf{p}_{t-1} \circ \mathbf{y}_t^j \in \mathcal{P}_t$.
- ¹⁶¹ **CLASSIFYING**: We calculate the score s_t^j of the image patches $\mathbf{x}_t^{\mathbf{p}'_t}$ corresponding to ¹⁶² all samples, or bounding boxes, using the current classifier θ_t ($h : \mathcal{X} \to \mathbb{R}$).

$$s_t^j = h(\mathbf{x}_t^{\mathbf{p}_{t-1} \circ \mathbf{y}_t^j} | \theta_t)$$
(1)

LABELING: We determine label l_t^j of each sample *j* using the score of the sample. If the score is above a threshold τ , the sample is likely to be target match,

$$l_t^j = \operatorname{sign}(s_t^j - \tau) \tag{2}$$

¹⁶⁵ **ESTIMATING**: We determine the next target state \mathbf{p}_t typically by selecting the best ¹⁶⁶ \mathbf{p}_t^j that corresponds to the maximum score s_t^j , $\mathbf{p}_t = \mathbf{p}_{t-1} \circ \mathbf{y}_t^*$

$$\mathbf{y}_{t}^{*} = \left\{ \mathbf{y}_{t}^{j^{*}} | j^{*} = \operatorname*{argmax}_{j \in \{1, \dots, N\}} s_{t}^{j} \right\}$$
(3)

¹⁶⁷ UPDATING: Finally, we update the classifier by its own labeled data,

$$\theta_{t+1} = u(\theta_t, \mathcal{X}_t, \mathcal{L}_t) \tag{4}$$

¹⁶⁸ in which u(l) is the update function (e.g., budgeted SVM update [12]), and \mathcal{X}_t , \mathcal{L}_t are ¹⁶⁹ the set of input patches and output labels used as the training set of the discriminator.

170 2.2 Baseline System Implementation

To develop a baseline tracking-by-detection algorithm for this study, we use a robust 171 part-based detector for the CLASSIFYING process. This detector employs strong 172 low-level features based on histograms of oriented gradients (HOG) and uses a la-173 tent SVM to perform efficient matching for deformable part-based models (pictorial 174 structures) [60]. From each frame, we draw N samples from a Gaussian distribution 175 whose mean is the target's bounding box in the last frame (including its location 176 and size). The selected detector then outputs the classification score for each sample, 177 which is thresholded to obtain the sample's label. The highest classification score is 178 considered as the current target location (Figure 1). 179



Figure 1: A simple tracking-by-detection pipeline. After gathering some samples from the current frame, the tracker employs its detector to label the samples as positive (target) or negative (background). The target position is estimated using these labeled samples. The labels, in turn, are used to update the classifier for the next frame.

In the first frame, we generate $\alpha_1 N$ positive samples by perturbing the first annotated target patch by few pixels in location and size, select $\alpha_2 N$ negative samples from local neighborhood of the target, and select $\alpha_3 N$ negative samples from global background of the object in a regular grid ($\alpha_1 + \alpha_2 + \alpha_3 = 1$). These samples are used to train the SVM detector in the first frame. From the next frames, the labels are obtained by the detector itself, and the classifier is batch-trained with all of the samples collected so far.

There are several parameters in the system such as the parameters of sampling step (number of samples *N*, effective search radius Σ_{search}). These parameters were tuned using a simulated annealing optimization on a cross-validation set. The partbase detector dictionary, and the thresholds τ_l , τ_u , and the rest of above-mentioned parameters have been adjusted using cross-validation. With N = 1000, $\tau = 0.34$ T1 achieved the speed of 47.29 fps on a Pentium IV PC @ 3.5 GHz and a Matlab/C++ implementation on a CPU.

¹⁹⁴ 2.3 Method of Evaluation

The experiments are conducted on 100 challenging video sequences, OTB-100 [61], which involves many visual tracking challenges such as illumination variation (IV), scale variation (SV), occlusions (OCC), deformations (DEF), motion blur (MB), fast motion (FM), in-plane rotation (IPR), out-of-play rotation (OPR), out-of-view problem (OV), background clutter (BC) and low resolution (LR). The performance of the trackers is compared with the area under the curve of success plots and precision plots, on all of the sequences, or a subset of them with the given attribute.

Success plot indicates the reliability of the tracker and its overall performance while precision plot reflects the accuracy of the localization. The area under the surface of this plot (*AUC*), counts the number of successes of tracker over time $t \in$ {1,..., *T*}, i.e. when the overlap of the tracker target estimation \mathbf{p}_t with the ground truth \mathbf{p}_t^* exceeds the threshold τ_{ov} . Success plot, graphs the success of the tracker



Figure 2: Quantitative performance comparison of the baseline tracker (T1), its variant without model update (T0), and the state-of-the-art trackers using success plot.

against different values of the threshold τ_{ov} and its AUC is calculated as

$$AUC = \frac{1}{T} \int_0^1 \sum_{t=1}^T \mathbb{1} \left(\frac{|\mathbf{p}_t \cap \mathbf{p}_t^*|}{|\mathbf{p}_t \cup \mathbf{p}_t^*|} > \tau_{ov} \right) d_{\tau_{ov}}, \tag{5}$$

where *T* is the length of sequence, |.| denotes the area of the region, \cap and \cup stands for intersection and union of the regions respectively, and $\mathbb{1}(.)$ denotes the step function that returns 1 iff its argument is positive and 0 otherwise. This plot provides an overall performance of the tracker, reflecting target loss, scale mismatches, and localization accuracy.

To establish a fair comparison with the state-of-the-art of tracking-by-detection algorithms, TLD [54] and STRUCK [12] are selected based on the results of [23], BSBT [62] and MIL [43] is selected based on popularity, and CSK [32] was selected as one of the latest algorithms in the category. Since our trackers contain random elements (in sampling and re-sampling), the results reported here is the average of five independent runs.

219 2.4 Results

Figure 2 presents the success and precision plots of T1 along with other competitive trackers for all sequences. We also included a fixed version of T1 tracker (a detector without model update) as T0 to emphasize the role of updating. The figure demonstrates that without the model update, the detector cannot reflect the changes in target appearance and lose the target rapidly in most of the scenarios (comparing T0 and T1). However, it is also evident that having a single tracker is not robust against all of the target's variations (in line with [56]) and the performance of T1 is still low.

227 3 Co-tracking

A single detector may have difficulties in distinguishing the target from the background in certain scenarios. In those cases, it is beneficial to consult another detector with higher robustness. These second detector may have complimentary characteristics to the first one, or simply may be a more sophisticated detector that trades computational complexity with speed.

Collaborative discriminative trackers utilize classifiers that exchange their information, to achieve more robust tracking. These information exchanges are in the form of queries that one classifier sends to another. The purpose of this information exchange is to bridge across long-term and short-term memories [58], accommodate multi-memory dictionaries [63], mixture of deep and shallow models [64], facilitate multi-view on the data [30], and enable learning from mistakes [54].

239 3.1 Formalization

Built on co-training principle [65], collaborative tracking (co-tracking) provides a framework in which two classifiers exchange their information to promote tracking results and break self-learning loop. In this two-classifier framework [30], the challenging samples for one classifier are labeled by the other one, i.e., if a classifier finds a sample difficult to label, it relies on the other classifier to label it for this frame and similar samples in the future. In this case, we calculate the discrimination score s_t^j as a weighted sum of the two discriminant functions, $s_t^j = \sum_{c=1}^2 \alpha_t^{(c)} h(\mathbf{x}_t^j | \theta_t^{(c)})$ where $\alpha_t^{(c)}$ denotes the weight of each discriminator $\theta_t^{(c)}$, c = 1, 2. At the **CLASSIFYING** step, the corresponding sample \mathbf{x}_t^j is considered as a challenging sample for the *c*-th discriminator when $\tau_l < h(\mathbf{x}_t^j | \theta_t^{(c)}) < \tau_u$ holds because it locates close to the corresponding discrimination boundary. When one of the two discriminators answered



Figure 3: Collaborative tracking. A detector and an auxiliary classifier trust each other to handle the sample difficult for them to classify.

it challenging, the score of the sample is calculated with using the other score.

$$s_{t}^{j} = \begin{cases} \alpha_{t}^{(2)} h(\mathbf{x}_{t}^{j} | \theta_{t}^{(2)}) &, \tau_{l} < h(\mathbf{x}_{t}^{j} | \theta_{t}^{(1)}) < \tau_{u} \\ \alpha_{t}^{(1)} h(\mathbf{x}_{t}^{j} | \theta_{t}^{(1)}) &, \tau_{l} < h(\mathbf{x}_{t}^{j} | \theta_{t}^{(2)}) < \tau_{u} \\ \sum_{c=1}^{2} \alpha_{t}^{(c)} h(\mathbf{x}_{t}^{j} | \theta_{t}^{(c)}) &, \text{otherwise} \end{cases}$$
(6)

At the **UPDATING** step, the weight $\alpha_t^{(c)}$ of the discriminator *c* is adjusted according to the degree of contradiction to the provisional answers that are determined at the **ESTIMATION** step by an integration of all the information. Finally, the classifiers are updated using only the samples that they successfully labeled in the previous frame to reflect the latest target changes.

245 3.2 Evaluation

For this experiment, we selected a naive classifier with complementary properties 246 to the main classifier in the previous section. This classifier is a KNN classifier us-247 ing HOC and HOG features, trained on the samples trained from the first frame, 248 and updated with all the labeled samples by the collaboration of the classifiers. Not 249 being pre-trained, the performance of this auxiliary classifier is poor in the begin-250 ning, but gradually gets better. The quick classification of the KNN (owning to its 251 kd-tree implementations and lightweight features) and lack of pre-training, grants it 252 high speed and generalization which is in contrast to the main detector. However, 253 it should be noted that without being supervised by the main SVM-based detector, 254 this classifier cannot perform well in isolation for tracking task. Figure ?? presents 255 the performance of this auxiliary tracker as T2. As observed in the figure, the perfor-256 mance of the obtained co-tracker (T3), is better than the main detector (T1) and the 257 auxiliary classifier (T2) as a result of co-labeling, data exchange, and co-learning. 258

²⁵⁹ 4 Active Co-tracking

The co-tracking framework provides a means for classifiers to exchange information. 260 This framework utilizes a utility measure (e.g., the classification confidence in [30]) 261 to select the data for which one of the collaborates fails to classify with high confi-262 dence, and then train the other classifier on those data. This approach has two main 263 shortcomings: (1) the redundant labeling of all samples for both classifiers and (2) 264 training the collaborator with "all" of the uncertain samples. While the former in-265 crease the complexity of the system, the latter is not the optimal solution for tracking 266 a target with non-stationary appearance distributions [31]. 267

In this view, a principled ordering of samples for training [66], and selecting a subset of them based on criteria [33] can reduce the cost of labeling leading to faster performance increase as a function of the amount of data available. It is found that

detectors trained with an effective, noise-free, and outlier-free subset of the training data may achieve higher performance than those trained with the full set [67, 68].

Robust learning algorithms provide an alternative way of differentially treating 273 training examples, by assigning different weights to different training examples or 274 by learning to ignore outliers [69]. Learning first from easy examples [70], pruning 275 adversarial examples¹ [71], and sorting the samples based on their training value 276 [33] are some of the approaches explored in the literature. However, the most com-277 mon setting is active learning, whereby most of the data is unlabeled and an algo-278 rithm selects which training examples to label at each step, for the highest gains in 279 performance. Thus, some active learning approaches focus on learning the hard-280 est examples first (those closest to the decision boundary). Some approaches focus 281 on learning the hardest examples first (e.g., those closest to the decision boundary), 282 whereas some others gauge the information contained in the sample and select the 283 most informative ones first. For example, Lewis and Gale [72] utilized the uncer-284 tainty of the classifier for a sample as an index of its usefulness for training. 285

²⁸⁶ **4.1** The idea

Active learning has been used in visual tracking to consider the uncertainty caused by bags of samples [51], to reduce the number of necessary labeled samples [73], to unify sample learning and feature selection procedure [74], and to reduce the sampling bias by controlling the variance [75].

In this study, we utilized the sampling uncertainty that can bind the active learn-291 ing and co-tracking. As mentioned earlier, the baseline classifier, despite being accu-292 rate, has low generalization on new samples, slow classification speed, and computa-293 tionally expensive re-training. On the other hand, the auxiliary classifier is agile and 294 learns rapidly, with negligible retraining time. To combine the merits of these two 295 classifiers, to cancel out their demerits with one another, and to address the afore-296 mentioned issues of co-tracking (redundant labeling and excessive samples), we in-297 corporate an active learning module to select the most informative data, i.e. those 298 for which the naive classifier is most uncertain, and query their labels from the part-299 based detector. This architecture (Figure 4, here called T4) mainly use naive classifier 300 for labeling the data and only ask the label of hard samples from the slower detec-301 tor, therefore, limits the redundancy and unleash the speed of the agile classifier. In 302 addition, by training the naive classifier only on hard samples, the generalization of 303 this classifier is preserved while increasing its accuracy. 304

To further increase the accuracy of the tracker and make it more robust against occlusions and drastic temporal changes of the target, it is possible to update the detector less frequently. This asymmetric version of the active co-tracker (T5), by introducing long-term memory to the tracker, benefits from combining the long and

¹Images with tiny, imperceptible perturbations that fool a classifier into predicting the wrong labels with high confidence.

input : Target position in last frame \mathbf{p}_{t-1} **output**: Target position in current frame \mathbf{p}_t

for $j \leftarrow 1$ to n do Generate a sample $\mathbf{p}_t^j \sim \mathcal{N}(\mathbf{p}_{t-1}, \Sigma_{search})$ Calculate $s_t^j \leftarrow h(\mathbf{x}_t^{\mathbf{p}_t^j} | \theta_t^{(1)})$ (eq(7)) Determine uncertain samples \mathcal{U}_t (eq(8)) if $\mathbf{p}_t^j \in \mathcal{U}_t$ then $\theta_t^{(1)}$ is uncertain $\left| \begin{array}{c} Query \ \theta_t^{(2)} : \ l_t^j \leftarrow Sign\left(h(\mathbf{x}_t^{\mathbf{p}_t^j} | \theta_t^{(2)})\right) \right)$ else $\left| \begin{array}{c} Label using \ \theta_t^{(1)} : \ l_t^j \leftarrow Sign(s_t^j) \\ \mathcal{D}_t \leftarrow \mathcal{D}_t \cup \langle \mathbf{x}_t^{\mathbf{p}_t^j}, \ l_t^j \rangle$ Update $\theta_t^{(2)}$ with $\mathcal{D}_{t-\Delta,..,t}$ every Δ frames ($\Delta = 1$ for T4) if $\sum_{j=1}^n \mathbbm{1}(l_t^j > 0) > \tau_p$ and $\sum_{j=1}^n \pi_t^j > \tau_a$ then $Approximate target state \ \mathbf{\hat{p}}_t$ (eq(10)) Update $\theta_t^{(1)}$ with \mathcal{U}_t else target occluded $| \ \mathbf{\hat{p}}_t \leftarrow \mathbf{p}_{t-1}$ Algorithm 1: Active Co-Tracking (ACT)

short-term collaboration (as in [58]) and reduces the frequency of the expensive up dates of the tracker (Algorithm 1).

311 4.2 Formalization

In the proposed active co-tracking framework, a main classifier attempts to label the sample, and it queries the label from the other classifier if the main classifier emit uncertain results. This is in contrast with using a linear combination of both classifiers based on their classification accuracy as adopted in T3. At the **CLASSIFYING** step, the proposed tracker can score each sample based on the classifier confidence, i.e., for sample \mathbf{p}_{t}^{j} we calculate score s_{t}^{j}

$$s_t^j = h(\mathbf{x}_t^{\mathbf{p}_t^j} | \boldsymbol{\theta}_t^{(1)}). \tag{7}$$

Based on uncertainty sampling [72], the samples for which the classification score is more uncertain (i.e., $s_t^j \rightarrow 0$), contains more information for the classifier if they are labeled by the other classifier. Therefore, the scores of all samples are sorted, and *m* samples with the closest values to 0 are selected to be queried from $\theta_t^{(2)}$. To handle the situations for which the number of highly uncertain samples are more than *m*,



Figure 4: Active co-tracker, a collaborative tracker that utilizes an active query mechanism to query the most informative samples from the main detector, and feed them to the light-weight classifier to learn.

³²³ a range of scores are determined by lower and higher thresholds (τ_l and τ_u) and all ³²⁴ the samples in this range are considered highly uncertain.

$$\mathcal{U}_t = \{ \mathbf{p}_t^i | \tau_l < s_t^i < \tau_u \text{ or } \mid \{ \exists j \neq i | s_t^j \le s_t^i \} \mid < m \}$$

$$\tag{8}$$

in which U_t is the list of uncertain samples. The label of the samples $l_t^j \in \mathcal{L}_t, j = 1, ..., N$ are then determined by

$$l_t^j = \begin{cases} sign\left(h\left(\mathbf{x}_t^{\mathbf{p}_t^j}|\theta_t^{(1)}\right)\right) &, \mathbf{p}_t^j \in \mathcal{U}_t \\ sign\left(h\left(\mathbf{x}_t^{\mathbf{p}_t^j}|\theta_t^{(2)}\right)\right) &, \mathbf{p}_t^j \notin \mathcal{U}_t \end{cases}$$
(9)

and all image patches $\mathbf{x}_{t}^{\mathbf{p}_{t}^{\prime}}$ and labels l_{t}^{j} are stored in \mathcal{D}_{t} .

At the **ESTIMATION** step, we follow the importance sampling mechanism originally employed by particle filter trackers,

$$\hat{\mathbf{p}}_{t} = \frac{\sum_{j=1}^{n} \pi_{t}^{j} \mathbf{p}_{t}^{j}}{\sum_{j=1}^{1} \pi_{t}^{j}}.$$
(10)

where $\pi_t^j = s_t^j \mathbb{1}(l_t^j > 0)$ and $\mathbb{1}(.)$ is the indicator function, 1 if true, zero other-328 wise. This mechanism approximates the state of the target, based on the effect of 329 positive samples, in which samples with higher scores gravitates the final results 330 more toward themselves. Upon the events such as massive occlusion or target loss, 331 this sampling mechanism degenerates [10]. In such cases, the number of positive 332 samples and their corresponding weights shrinks significantly, and the importance 333 sampling is prone to outliers, distractors, and occluded patches. To address this is-334 sue, if the number of positive samples is less than τ_p , and their score average is less 335 than τ_a , the target is deemed occluded to avoid tracker degeneracy. 336



Figure 5: Quantitative performance comparison of the asymmetric active co-tracker (T5), active co-tracker (T4), the ordinary co-tracker (T3), and their individual trackers (T1 and T2).

337 4.3 Evaluation

Figure 5 illustrates the effectiveness of the proposed trackers against their baselines. The active query mechanism in T4 improves the efficiency and effectiveness of cotracking (T3). Especially in the asymmetric co-tracker (T5), the mixture of long-term and short-term memory classifiers using this method is to key to automatically balance the stability-plasticity equilibrium. It is also prudent for the tracker to adapt to the temporal distribution of the target appearance, before its re-distribution by illumination changes, etc.

In summary, the advantage of the proposed trackers especially the asymmetric 345 ones (T5) compared to the conventional co-tracking (T3) are as follows: (1) the clas-346 sifiers do not exchange all the data they have problems in labeling, instead, the most 347 informative samples are selected by uncertainty sampling, and exchanged. (2) the 348 update rate of classifiers is different to realize a short and long-term memory mix-349 ture, (3) the samples that are labeled for the target localization can be re-used for 350 training and the need for an extra round of sampling and labeling is revoked, (4) 351 since in the proposed asymmetric co-tracking, one of the classifiers scaffolds the 352 other one instead of participating in every labeling process, a more sophisticated 353 classifier with higher computational complexity can be used. 354

5 Active Ensemble Co-tracking

Ensemble discriminative tracking utilizes a committee of classifiers, to label data
samples, which are in turn, used for retraining the tracker to localize the target using
the collective knowledge of the committee. In such frameworks the labeling process
is performed by leveraging a group of classifiers with different views [41, 76, 52],
subsets of training data [37, 77] or memories [53, 78].

In ensemble tracking [56, 41, 79, 80, 43, 52, 81, 53] the self-learning loop is broken, and the labeling process is performed by eliciting the belief of a group of classifiers. However, this framework typically do not address some of the demands of tracking-by-detection approaches like a proper model update to avoid model drift or non-stationary of the target sample distribution. Besides, ensemble classifiers do not exchange information, and collaborative classifiers entirely trust the other classifier to label the challenging samples for them and are susceptible to label noise.

Traditionally, ensemble trackers were used to providing a multi-view classifica-368 tion of the target, realized by using different features to construct weak classifiers. 369 In this view, different classifiers represent different hypotheses in the version-space, 370 to accurately model the target appearance. Such hypotheses are highly-overlapping, 371 therefore an ensemble of them overfits the target. The desired committee, however, 372 consists of competing hypotheses, all consistent with the training data, but each of 373 the specialized in certain aspect. In this view, the most informative data samples 374 are those about which the hypotheses disagree the most, and by labeling them the 375 version-space is minimized leading to quick convergence yet accurate classification 376 [82]. Motivated by this, we proposed a tracker that employs a randomized ensemble 377 of classifiers and selects the most informative data samples to be labeled. 378

379 5.1 The idea

One of the most theoretically-motivated query selection frameworks is Query-by-Committee (QBC) algorithm [82, 83] that maintain a committee of models which are all trained on the current labeled set, but represent competing hypotheses. Each committee member is then allowed to vote on the labeling of query candidates. The most informative query is considered to be the instance about which they most disagree. The premise behind the QBC framework is minimizing the size of the *version space*, which is the set of hypotheses that are consistent with data.

The original QBC was built upon randomized component learning algorithm. For other model classes, such as discriminative or non-probabilistic models, Abe and 388 Mamitsuka [84] have proposed Query-by-Bagging (QBag), which employ bagging 389 [85] to construct committees. Bagging is a technique to enhance the performance 390 of the existing learning algorithm by running it many times on a set of re-sampled 391 governed by a uniform distribution and the final hypothesis is obtained by taking 392 majority vote over the output of predictions of the output hypotheses. QBag intro-393 duces the randomness in the form of re-sampling from the input data based on the 394 idea that prediction error consists of *bias*, which is the estimation error due to the 395 smaller input size, and *variance* which is explained by the statistical variation exist-396 ing in data. Bagging can isolate bias from variance and minimize the latter [84]. 397

We propose the adjustment of the QBag algorithm for online training to solve the label noise problem in T6. Similar to T5, the drift problem is handled using dualmemory strategy: the committee rapidly adapts to target changes, whereas the main classifier possesses a longer memory to promote the stability of the target template.



Figure 6: Active ensemble co-tracker. The bagging-induced ensemble labels the input samples and only query the most disputed ones from the slow part-based classifier.

402 **5.2 Formalization**

⁴⁰³ An ensemble discriminative tracker employs a set of classifiers instead of one. These ⁴⁰⁴ classifiers, hereafter called *committee*, are represented by $C = \{\theta_t^{(1)}, \ldots, \theta_t^{(C)}\}$, and are ⁴⁰⁵ typically homogeneous and independent (e.g., [52, 81]). Popular ensemble trackers

⁴⁰⁶ utilize the majority voting of the committee as their utility function,

$$\mathbf{s}_t^j = \sum_{c=1}^C \operatorname{sign}\left(h(\mathbf{x}_t^{\mathbf{p}_{t-1} \circ \mathbf{y}_t^j} | \boldsymbol{\theta}_t^{(c)})\right). \tag{11}$$

and eq(9) is used to label the samples. Finally, the model is updated for each classifier independently, meaning that each of the committee members are trained with a random subset of the uncertain set. $\theta_{t+1}^{(c)} = u(\theta_t^{(c)}, \Gamma_t^{(c)} \sim U_t)$ where $u(\theta, \mathcal{X})$ is the updating the model θ with samples \mathcal{X} . The uncertain set U_t contains all of the samples for which the ensemble disagree and were sent to the auxiliary classifier for labeling. The detector $\theta_t^{(o)}$ is also updated with all recent data $\mathcal{D}_{t-\Delta,..,t}$ every Δ frames.

413 5.3 Evaluation

Figure 7 depicts the overall performance of the proposed tracker against other benchmarked algorithms on all sequences of the dataset. The plots show that T6 has a superior performance over T5 and its predecessors. It also reveals that the tracker has many accurate estimations of the target (sharp slope between $0.9 \ge \tau_{ov} > 1$). Furthermore, the other steep slope around $\tau_{ov} \approx 0.4$ and the high value when $\tau_{ov} \rightarrow 0$ suggest that tracker was able to keep track of the target in most cases, and the devised scheme effectively reduced the drift problem.

421 6 Discussion

The instances of the proposed framework are evaluated against state-of-the-art trackers on public sequences that become the de-facto standards of benchmarking the



Figure 7: Quantitative performance comparison of the active ensemble co-tracker (T6) with its predecessors.

trackers. The trackers are compared with popular metrics such as success plot and 424 precision plot to establish a fair benchmark. In addition, the performance of the pro-425 posed trackers are investigated for videos with a distinguished tracking challenge, 426 and the results are compared with state-of-the-art and discussed. Additionally, the 427 effect of the information exchanged will be examined thoroughly to illustrate the dy-428 namics of the system. The preliminary results of the proposed framework demon-429 strate a superior performance for the proposed trackers when applied on all the se-430 quences, and most of the subsets of the test dataset with distinguished challenges. 431 Finally, the future research direction is discussed and the opened research avenues 432 are introduced to the field. 433 As Figure 7 and Table 2 demonstrates, T6 has the best overall performance among 434 investigated trackers on this dataset. While this algorithm has a clear edge in han-435

⁴³⁵ Investigated trackers on this dataset. While this algorithm has a clear edge in han ⁴³⁶ dling many challenges, its performance is comparable with T5 in the case of occlu ⁴³⁷ sions and z-rotations. It is also evident that T6 is troubled with fast deformations
 ⁴³⁸ since neither of the ensemble members is specialized in handling a specific type of
 ⁴³⁹ deformations and the collective decision of the ensemble may involve mistakes with

high confidence. On the other hand, T5 utilizes a dual memory scheme and a single 440 classifier can handle extreme temporal deformations better than the ensemble in T6. 441 Interestingly, it is observed that in most of the subcategories that T6 is clearly better 442 than the other trackers, the success plot of T6 starts with a plateau and later has a 443 sharp drop around $\tau_{ov} = 0.8$. This means that T6 provides high-quality localization 444 (i.e., bigger overlaps with the ground truth). Similarly, from precision plot, it is ev-445 ident that T6 shows a graceful degradation in different scenarios, and although it 446 does not provide a good scale adaptation for targets, it is able to localize them better 447 than the competing trackers. 448

Table 2: Quantitative evaluation of state-of-the-art under different visual tracking challenges using AUC of success plot (%). The first, second and third best methods are shown in color.

	IV	DEF	OCC	SV	IPR	OPR	ov	LR	BC	FM	MB	ALL
Т0	12	12	13	12	13	13	14	5	12	15	18	14
T1	37	29	3	36	42	39	43	30	33	39	36	38
T2	23	19	23	23	28	25	25	22	23	24	20	25
T3	41	32	39	40	44	42	43	30	36	43	39	41
T4	50	39	47	48	53	49	48	37	44	50	45	49
T5	52	47	53	51	59	56	52	38	41	53	46	52
T6	57	40	51	53	61	55	63	46	53	60	58	56
TLD	49	32	42	44	50	43	45	37	40	45	42	46
STRK	46	41	44	43	51	48	44	39	39	52	48	48
CSK	40	36	36	34	43	39	32	29	42	39	32	41
MIL	35	35	38	35	41	39	40	32	31	35	28	36
BSBT	23	18	23	21	27	24	32	23	23	26	24	25

449 7 Conclusions and Future Works

This chapter provides a step-by-step tutorial for creating an accurate and high-performance 450 tracking-by-detection algorithm out of ordinary detectors, by eliciting an effective 451 collaboration among them. The use of active learning in junction with co-learning 452 enable the creation of a battery of tracker that strive to minimize the uncertainty of 453 one classifier by the help of another. Finally, we proposed to employ a committee of 454 classifiers, each trained incrementally on a randomized portion of the latest obtained 455 training samples, to enhance the discriminative power of the tracker. This idea is in-456 spired by the query-by-bagging framework that follows the version-space shrinking 457 strategy to distinguish the most informative samples. Such samples are then queried 458 from a more complicated classifier with longer memory that is robust against fluctu-459 ations in target appearance and occlusions. Furthermore, using an expectation of the 460 bounding boxes compensates for over-reliance of the tracker on the classifiers con-461 fidence function. The balance in stability-plasticity equilibrium is achieved by the 462 combination of several short-term classifiers with a long-term classifier, and manag-463 ing their interaction with an active learning mechanism. 464



(a) Tracking results of sequence FaceOcc2 and Walking2 with severe occlusions



(b) Tracking results of sequence Deer and Jumping with motion blur



(c) Tracking results of sequence Girl and Ironman with in-plane and out-of-plane rotations



(d) Tracking results of sequence Singer1, CarDark and Shaking with drastic illumination changes



(e) Tracking results of sequence Board with background clutter

Figure 8: Sample tracking results of evaluated algorithms on several challenging video sequences. In these sequences the red box depicts the T6 against other trackers (T0-5 in blue and TLD, STRK, CSK, MIL, and BSBT in gray). The ground truth is illustrated with yellow dashed box. The results are available in the http://ishiilab.jp/member/meshgi-k/act.html.

The trail of proposed trackers led to T6, which incorporates ensemble tracking,
 active learning, and co-learning in a discriminative tracking framework and outper form state-of-the-art discriminative and generative trackers on a large video dataset
 with various types of challenges such as appearance changes and occlusions.

The future direction of this study involves including other detectors to care for context, to have accurate physical models for known categories, to use deep features to improve discrimination and to examine different methods of building the ensemble and detecting most informative samples or exchanging.

473 Acknowledgment

⁴⁷⁴ This article is based on results obtained from a project commissioned by the Japan

⁴⁷⁵ NEDO and was supported by Post-K application development for exploratory chal-



⁴⁷⁶ lenges from the Japan MEXT.

477 **References**

- [1] J. Cech, R. Mittal, A. Deleforge, and R. Horaud, "Active-speaker detection and localization with mic and cameras embedded into a robotic head," in *Humanoids*'13, 2013. 1
- [2] A. Cosgun, D. A. Florencio, and H. I. Christensen, "Autonomous person following for telepresence robots," in *ICRA'13*. IEEE, 2013, pp. 4335–4342.
- [3] T. B. Moeslund, A. Hilton, and V. Krüger, "A survey of advances in vision-based human motion capture and analysis," *CVIU*, vol. 104, no. 2, pp. 90–126, 2006. 1
- [4] M. Störring, T. B. Moeslund, Y. Liu, and E. Granum, "Computer vision-based gesture
 recognition for an augmented reality interface," in *VIIP'04*, vol. 3, 2004, pp. 766–771.
- [5] O. Koller, O. Zargaran, H. Ney, and R. Bowden, "Deep sign: Hybrid cnn-hmm for continuous sign language recognition," in *BMVC'16*, 2016. 1
- [6] L. Wang, B. Schmidt, and A. Y. Nee, "Vision-guided active collision avoidance for human robot collaborations," *Manufacturing Letters*, vol. 1, no. 1, pp. 5–8, 2013.
- [7] A. Ess, B. Leibe, K. Schindler, and L. Van Gool, "Moving obstacle detection in highly
 dynamic scenes," in *ICRA'09*. IEEE, 2009, pp. 56–63. 1
- [8] L. Xia, I. Gori, J. K. Aggarwal, and M. S. Ryoo, "Robot-centric activity recognition from
 first-person rgb-d videos," in WACV'15. IEEE, 2015, pp. 357–364. 1
- [9] S. S. Rautaray and A. Agrawal, "Vision based hand gesture recognition for human com puter interaction: a survey," *AI Review*, vol. 43, no. 1, pp. 1–54, 2015. 1
- [10] C. Bao, Y. Wu, H. Ling, and H. Ji, "Real time robust l1 tracker using accelerated proximal
 gradient approach," in *CVPR'12*. 1, 2, 12
- [11] J. Kwon and K. M. Lee, "Tracking by sampling trackers," in ICCV'11, 2011. 1
- [12] S. Hare, A. Saffari, and P. H. Torr, "Struck: Structured output tracking with kernels," in
 ICCV'11, 2011. 1, 2, 4, 5, 7
- [13] T. B. Dinh, N. Vo, and G. Medioni, "Context tracker: Exploring supporters and distracters
 in unconstrained environments," in *CVPR'11*, 2011. 1, 3
- [14] K. Meshgi, S.-I. Maeda, S. Oba, and S. Ishii, "Data-driven probabilistic occlusion mask to
 promote visual tracking," in *CRV'16*, 2016. 1, 3
- [15] Y. Wu, H. Ling, J. Yu, F. Li, X. Mei, and E. Cheng, "Blurred target tracking by blur-driven tracker," in *ICCV*'2011, 2011.
- ⁵⁰⁷ [16] D. A. Ross, J. Lim, R.-S. Lin, and M.-H. Yang, "Incremental learning for robust visual ⁵⁰⁸ tracking," *IJCV*, vol. 77, no. 1-3, pp. 125–141, 2008. 1
- [17] J. Fang, H. Xu, Q. Wang, and T. Wu, "Online hash tracking with spatio-temporal saliency auxiliary," *CVIU*, 2017. 1, 3
- [18] A. Taalimi, H. Qi, and R. Khorsandi, "Online multi-modal task-driven dictionary learn ing and robust joint sparse representation for visual tracking," in AVSS'15, 2015. 1
- [19] H. Kiani, T. Sim, and S. Lucey, "Correlation filters with limited boundaries," in CVPR'15, 2015. 2
 - 19

- [20] M. Danelljan, G. Hager, F. Shahbaz Khan, and M. Felsberg, "Learning spatially regularized correlation filters for visual tracking," in *ICCV'15*, 2015, pp. 4310–4318.
- [21] M. Danelljan, A. Robinson, F. S. Khan, and M. Felsberg, "Beyond correlation filters: Learning continuous convolution operators for visual tracking," in ECCV'16. 2
- [22] H. Nam and B. Han, "Learning multi-domain convolutional neural networks for visual tracking," in CVPR'16. 2
- [23] Y. Wu, J. Lim, and M.-H. Yang, "Online object tracking: A benchmark," in *CVPR'13*.
 IEEE, 2013, pp. 2411–2418. 2, 7
- [24] M. Kristan, J. Matas, A. Leonardis, and M. Felsberg, "The visual object tracking vot2015
 challenge results," in *ICCVw'15*. 2
- [25] A. Li, M. Lin, Y. Wu, M.-H. Yang, and S. Yan, "Nus-pro: A new visual tracking challenge,"
 PAMI, 2016. 2
- [26] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in CVPR'16, 2016, pp. 770–778. 2
- [27] N. Wang and D.-Y. Yeung, "Learning a deep compact image representation for visual tracking," in *NIPS'13*, 2013, pp. 809–817.
- [28] H. Li, Y. Li, F. Porikli *et al.*, "Deeptrack: Learning discriminative feature representations
 by convolutional neural networks for visual tracking." in *BMVC*, vol. 2014, 2014. 2
- [29] S. Hong, T. You, S. Kwak, and B. Han, "Online tracking by learning discriminative saliency map with convolutional neural network," in *ICML'15*, 2015, pp. 597–606. 2
- [30] F. Tang, S. Brennan, Q. Zhao, and H. Tao, "Co-tracking using semi-supervised support vector machines," in *ICCV'07*, 2007. 2, 4, 8, 9
- [31] Q. Bai, Z. Wu, S. Sclaroff, M. Betke, and C. Monnier, "Randomized ensemble tracking," in *ICCV*'13, 2013. 2, 9
- [32] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista, "Exploiting the circulant structure of tracking-by-detection with kernels," in *ECCV'12*. Springer, 2012, pp. 702–715. 2, 7
- [33] A. Lapedriza, H. Pirsiavash, Z. Bylinskii, and A. Torralba, "Are all training examples
 equally valuable?" arXiv, 2013. 2, 9, 10
- 543 [34] I. Matthews, T. Ishikawa, and S. Baker, "The template update problem," PAMI, 2004. 2
- [35] H. Grabner, C. Leistner, and H. Bischof, "Semi-supervised on-line boosting for robust tracking," in ECCV'08, 2008. 2, 3, 4
- [36] K. Meshgi, M. S. Mirzaei, S. Oba, and S. Ishii, "Efficient asymmetric co-tracking using uncertainty sampling," in *ICSIPA'17*, 2017. 3
- [37] K. Meshgi, S. Oba, and S. Ishii, "Robust discriminative tracking via query-bycommittee," in AVSS'16, 2016. 3, 4, 13
- [38] P. Pérez, C. Hue, J. Vermaak, and M. Gangnet, "Color-based probabilistic tracking," in
 ECCV'02. 3
- [39] Y. Wu, M. Pei, M. Yang, and Y. Jia, "Robust discriminative tracking via landmark-based
 label propagation," *TIP*, 2015. 3, 4
- [40] N. C. Oza and S. Russell, "Online ensemble learning," in AAAI'00. 3
 - 20

- [41] H. Grabner, M. Grabner, and H. Bischof, "Real-time tracking via on-line boosting." in
 BMVC'06, vol. 1, no. 5, 2006, p. 6. 3, 4, 13, 14
- ⁵⁵⁷ [42] C. Leistner, A. Saffari, P. Roth, and H. Bischof, "On robustness of on-line boosting: a ⁵⁵⁸ competitive study," in *ICCVw*'09. **3**
- [43] B. Babenko, M.-H. Yang, and S. Belongie, "Visual tracking with online multiple instance learning," in *CVPR'09*, 2009. 3, 7, 14
- ⁵⁶¹ [44] S. Avidan, "Support vector tracking," PAMI, vol. 26, no. 8, pp. 1064–1072, 2004. 3
- [45] H. Masnadi-Shirazi, V. Mahadevan, and N. Vasconcelos, "On the design of robust classi fiers for computer vision," in *CVPR'10*, 2010. 3
- [46] C. Leistner, A. Saffari, J. Santner, and H. Bischof, "Semi-supervised random forests," in
 ICCV'09. 3
- [47] B. Zeisl, C. Leistner, A. Saffari, and H. Bischof, "On-line semi-supervised multipleinstance boosting," in *CVPR'10.* 3
- [48] K. Zhang and H. Song, "Real-time visual tracking via online weighted multiple instance
 learning," *PR*, 2013. 3
- [49] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista, "High-speed tracking with kernelized correlation filters," *PAMI*, vol. 37, no. 3, pp. 583–596, 2015. 3
- ⁵⁷² [50] H. Grabner, J. Matas, L. Van Gool, and P. Cattin, "Tracking the invisible: Learning where
 ⁵⁷³ the object might be," in *CVPR'10*. IEEE, 2010, pp. 1285–1292. 3
- [51] K. Zhang, L. Zhang, M.-H. Yang, and Q. Hu, "Robust object tracking via active feature selection," *IEEE CSVT*, vol. 23, no. 11, pp. 1957–1967, 2013. 3, 10
- [52] A. Saffari, C. Leistner, J. Santner, M. Godec, and H. Bischof, "On-line random forests," in ICCVw'09. 4, 13, 14, 15
- [53] J. Zhang, S. Ma, and S. Sclaroff, "Meem: Robust tracking via multiple experts using en tropy minimization," in ECCV'14. 4, 13, 14
- [54] Z. Kalal, K. Mikolajczyk, and J. Matas, "Tracking-learning-detection," *PAMI*, vol. 34, no. 7, pp. 1409–1422, 2012. 4, 7, 8
- [55] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *CVPR'14*, 2014, pp. 580–587.
- ⁵⁸⁴ [56] S. Avidan, "Ensemble tracking," PAMI, vol. 29, 2007. 4, 7, 14
- [57] T. Woodley, B. Stenger, and R. Cipolla, "Tracking using online feature selection and a local generative model." in *BMVC'07*, 2007. 4
- [58] Z. Hong, Z. Chen, C. Wang, D. Prokhorov, and D. Tao, "Multi-store tracker (muster): a
 cognitive psychology inspired approach to object tracking," in *CVPR'15*, 2015. 4, 8, 11
- [59] J. Li, Z. Hong, and B. Zhao, "Robust visual tracking by exploiting the historical tracker
 snapshots," in *ICCVW'15*, 2015, pp. 41–49. 4
- [60] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with
 discriminatively trained part-based models," *PAMI*, vol. 32, 2010. 5
- [61] Y. Wu, J. Lim, and M.-H. Yang, "Object tracking benchmark," PAMI, 2015. 6
- [62] S. Stalder, H. Grabner, and L. Van Gool, "Beyond semi-supervised tracking: Tracking should be as simple as detection, but not simpler than recognition," in *ICCVw'09*. 7
 - 21

- [63] J. Xing, J. Gao, B. Li, W. Hu, and S. Yan, "Robust object tracking with online multi-lifespan dictionary learning," in *ICCV'13*, 2013, pp. 665–672.
- [64] B. Zhuang, L. Wang, and H. Lu, "Visual tracking via shallow and deep collaborative model," *Neurocomputing*, vol. 218, pp. 61–71, 2016.
- [65] A. Blum and T. Mitchell, "Combining labeled and unlabeled data with co-training," in
 COLT'98, 1998. 8
- [66] S. Vijayanarasimhan and K. Grauman, "Cost-sensitive active visual category learning,"
 IJCV, 2011. 9
- [67] N. Razavi, J. Gall, P. Kohli, and L. Van Gool, "Latent hough transform for object detection," ECCV'12. 10
- [68] X. Zhu, C. Vondrick, D. Ramanan, and C. C. Fowlkes, "Do we need more training data or better models for object detection?." in *BMVC*'12. 10
- [69] F. De la Torre and M. J. Black, "Robust principal component analysis for computer vision," in *ICCV'01*, 2001.
- [70] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *ICML'09*, 2009. 10
- [71] J. Lu, T. Issaranon, and D. Forsyth, "Safetynet: Detecting and rejecting adversarial exam ples robustly," *arXiv*, 2017. 10
- [72] D. D. Lewis and W. A. Gale, "A sequential algorithm for training text classifiers," in ACM
 SIGIR'94, 1994, pp. 3–12. 10, 11
- [73] C. H. Lampert and J. Peters, "Active structured learning for high-speed object detection,"
 in *PR*. Springer, 2009, pp. 221–231. 10
- [74] C. Li, X. Wang, W. Dong, J. Yan, Q. Liu, and H. Zha, "Active sample learning and feature selection: A unified approach," arXiv preprint arXiv:1503.01239, 2015. 10
- [75] A. Beygelzimer, S. Dasgupta, and J. Langford, "Importance weighted active learning," in
 ICML'09. ACM, 2009, pp. 49–56. 10
- [76] B. Han, J. Sim, and H. Adam, "Branchout: Regularization for online ensemble tracking
 with convolutional neural networks," in *ICCV*'17, 2017, pp. 2217–2224. 13
- [77] K. Meshgi, S. Oba, and S. Ishii, "Efficient version-space reduction for visual tracking,"
 CRV'17, 2017. 13
- [78] K. Meshgi, S. Oba, and S. Ishii, "Active discriminative tracking using collective memory,"
 in MVA'17. 13
- [79] N. C. Oza, "Online bagging and boosting," in SMC'05, 2005. 14
- [80] A. Saffari, C. Leistner, M. Godec, and H. Bischof, "Robust multi-view boosting with priors," in ECCV'10, 2010. 14
- [81] C. Leistner, A. Saffari, and H. Bischof, "Miforests: Multiple-instance learning with ran domized trees," in *ECCV'10*, 2010. 14, 15
- [82] S. Seung, M. Opper, and H. Sompolinsky, "Query by committee," in COLT'92, 1992. 14
- [83] B. Settles, Active learning. Morgan & Claypool Publishers, 2012. 14
- [84] N. Abe and H. Mamitsuka, "Query learning strategies using boosting and bagging," in
 ICML'98, 1998. 14
- [85] L. Breiman, "Bagging predictors," Machine learning, 1996. 14
 - 22