

From Handheld to Unconstrained Object Detection: a Weakly-supervised On-line Learning Approach

Elisa Maiettini^{*1}, Andrea Maracani^{*1,2,3} Raffaello Camoriano², Giulia Pasquale¹, Vadim Tikhanoff⁴, Lorenzo Rosasco^{2,3,5} and Lorenzo Natale¹

Abstract—Deep Learning (DL) based methods for object detection achieve remarkable performance at the cost of computationally expensive training and extensive data labeling. Robots embodiment can be exploited to mitigate this burden by acquiring automatically annotated training data via a natural interaction with a human showing the object of interest, handheld. However, learning solely from this data may introduce biases (the so-called domain shift), and prevents adaptation to novel tasks. While Weakly-supervised Learning (WSL) offers a well-established set of techniques to cope with these problems in general-purpose Computer Vision, its adoption in challenging robotic domains is still at a preliminary stage. In this work, we target the scenario of a robot trained in a teacher-learner setting to detect handheld objects. The aim is to improve detection performance in different settings by letting the robot explore the environment with a limited human labeling budget. We compare several techniques for WSL in detection pipelines to reduce model re-training costs without compromising accuracy, proposing solutions which target the considered robotic scenario. We show that the robot can improve adaptation to novel domains, either by interacting with a human teacher (Active Learning) or with an autonomous supervision (Semi-supervised Learning). We integrate our strategies into an on-line detection method, achieving efficient model update capabilities with few labels. We experimentally benchmark our method on challenging robotic object detection tasks under domain shift. Code will be released for reproducibility at camera-ready stage.

I. INTRODUCTION

In the state-of-the-art, object detection is typically addressed with DL-based approaches [1], [2] that achieve remarkable performance. Despite their high accuracy, they are constrained by requiring long training times and large annotated datasets, limiting their adoption in such applied settings where quick adaptation to novel tasks is required. In Robotics, the embodiment of a robotic agent can be exploited to interact with the environment, including humans, to mitigate this burden and actively acquire training data. Regarding the interaction with humans, past work shows that a teacher-learner scenario can be exploited to automatically collect labeled images for object recognition [3] and detection [4]. Specifically, in those works the human teacher shows an object, while holding it in their hand, to the robot and 3D information is used to automatically

collect the location information. However, while effective and allowing for a natural interaction, this approach supports limited generalization to novel, unseen, scenarios [4], [5]. A further possibility is to exploit robots ability to navigate and autonomously explore the environment, acquiring training images during operation. Such images come in streams and can carry useful information, eventually containing the objects of interest, but they are not labeled. *Weakly-supervised Learning* (WSL) [6], is a well-established general purpose Computer Vision framework which targets learning from partially-annotated datasets. However, despite initial work in robotic vision [7], [5] the robotic literature misses a thorough comparison that investigates advantages and limitations of existing techniques, especially in the context considered in this paper. For instance, in [5], the unlabeled images are processed with a pre-trained model to either select the hard ones and ask a human expert to help and annotate them (*Active Learning* (AL) framework [8]) or add the predictions of the easy ones to the training set (*Semi-supervised Learning* (SSL) [9], [6]). These frameworks allow for a natural interaction with the environment and the human teacher to improve the visual system and work presented in [5] effectively reduces the amount of manual annotation, but it has some limitations. Firstly, the unsupervised data processing is *pool-based* [8], that is, all unlabeled images are evaluated before query selection. This is not suitable for a robotic system that is exploring the environment and needs to decide interactively whether to request annotations or not. To this aim, *stream-based* techniques [8] are preferable, because they allow to process images frame by frame and to make individual query decisions on-line. This strategy, however, might yield to lower accuracy since queries are constructed using limited information on the unlabeled set [8]. Moreover, the pre-trained detection method in [5] iterates multiple times over the unlabeled data, which, while allowing to refine the data selection, slows down learning. Finally, while succeeding in reducing the human effort required for refinement, [5] still needs a relatively high number of manual annotations, which prevents its adoption in on-line applications.

In this paper, we study how WSL techniques can be used to exploit the robot interaction with the environment and the human teacher to update and improve performance of object detection models previously trained with data of handheld objects. We focus on the stream-based scenario with the aim of increasing the human labeling efficiency of weakly-supervised on-line object detection. Moreover, we consider the case in which only one pass over the unlabeled

*Equal contribution

¹Humanoid Sensing and Perception, Istituto Italiano di Tecnologia, Genoa, Italy

²IIT@MIT - Laboratory for Computational and Statistical Learning, IIT, Genoa, Italy, and MIT, Cambridge, MA, USA

³MaLGA & DIBRIS, Università degli Studi di Genova, Genoa, Italy

⁴iCub Tech, Istituto Italiano di Tecnologia, Genoa, Italy

⁵Center for Brains, Minds and Machines, MIT, Cambridge, MA, USA

data is allowed. The main contributions of this work are as follows. We present and empirically evaluate several AL techniques for detection, typically used in general purpose computer vision. We compare *pool-based* and *stream-based* AL in challenging robotic scenarios and propose a solution to overcome limitations of the latter. We also consider the case where no human labeling is allowed for adaptation. Specifically, we investigate the domain shift effects occurring when using a model trained on data of handheld objects in different settings and how wrongly self-annotated data can degrade accuracy in those cases. Finally, we propose an SSL sampling method to overcome this problem and we empirically demonstrate that, in case no labeling is allowed, it can effectively improve model performance. This paper is organized as follows: we introduce WSL in Sec. II and we cover related work in Sec. III. In Sec. IV, we present our efficient detection methods, which are analyzed and validated in Sec. V. Sec. VI concludes the paper.

II. BACKGROUND

The supervised learning approach to object detection is centered on learning the detector function from an annotated (supervised) dataset $S_n = \{(x_i, Y_i)\}_{i=1}^n$ of images (x_i) and corresponding bounding boxes and labels annotations (Y_i). The methods described in Sec. III-A fall in this category. They contributed to a clear progress in detection accuracy and prediction speed. However, they need expensively-annotated large-scale datasets to be optimized. This property does not meet the robotic requirement for a detector to adapt to a variety of tasks, potentially unknown a-priori, in a short time span. However, while large annotated datasets might not be available, plenty of unsupervised images are usually accessible to robots. In this context, a training set $S_n = L \cup U$ is typically composed of a labeled subset $L = \{(x_i, Y_i)\}_{i=1}^{n_L}$ and an unlabeled subset $U = \{x_i\}_{i=1}^{n_U}$. WSL allows the agent to select unsupervised images from U and acquire their labels semi-autonomously for updating the detector, minimizing human effort and improving accuracy. WSL includes several subclasses of methods, depending on the label-acquisition mechanism [9], [6]. The most relevant for this work are Active and Semi-supervised Learning.

Active Learning. AL [8] interactively queries unsupervised examples for expert labeling to minimize human annotation and maximize accuracy. Unlabeled examples are chosen from U according to a *scoring function* and a *sampling strategy*. Their labels are then queried to an expert, and newly-annotated examples are added to L for training. If all images in U are accessible at selection time, sampling is referred to as *pool-based*. Otherwise, if only one candidate from U is accessible, sampling becomes a binary decision on keeping or dropping it and is called *stream-based*. The AL selection criterion we focus on is *uncertainty sampling*, which picks the examples the model is *least confident* about.

Semi-supervised Learning. In SSL [6], unlabeled images are annotated by the detector itself with no human intervention, propagating predicted labels to high-confidence regions of the input space by exploiting the geometry of the input

data distribution. This technique is effective if the detector is not overconfident of its predictions and if the confidence threshold for propagating predicted labels is strict enough.

III. RELATED WORK

A. Object Detection

Early approaches to object detection were based on feature dictionaries [10] or specific kinds of image descriptors [11]. Feature vectors were separately classified by supervised learning methods. Despite yielding limited accuracy, these approaches had the advantage of being parsimonious in terms of computations and dataset size. More recently, object detection experienced significant progress thanks to the introduction of DL-based methods. This determined clear improvements in terms of predictive performance, mainly due to the powerful representation capabilities of deep networks. Such approaches include two-stage detectors based on Region Proposal Networks (RPNs) [12] (like e.g. Faster R-CNN [12] and Mask R-CNN [1]) and related extensions [13], [14], [15]. These methods employ a deep network to perform (i) region candidates predictions, (ii) per-region feature extraction and (iii) region classification and refinement. Alternative end-to-end approaches include one-stage detectors, which replace the RPN with a fixed, dense grid of candidate bounding boxes. One such example is SSD [16], [17], achieving accuracies competitive with the RPN-based Faster R-CNN and high frame rate. Another one-stage method, RetinaNet [18], rebalances foreground and background examples through the so-called Focal Loss.

B. Efficient Object Detection for Robotics

Despite their high accuracy, the approaches described above typically require (i) long training time and (ii) large-scale annotated datasets for adaptation to novel tasks. These aspects limit their adoption in Robotics.

Computational efficiency. A well-known issue of DL-based pipelines is that they suffer from *catastrophic forgetting* when optimized on new data [19]. This limitation implies retraining these models on the full dataset, causing long adaptation time. To address this issue, a recent work for robotic object detection leverages fast classifiers to enable on-line adaptation [20], [21]. Specifically, in [21], an efficient multi-stage pipeline is proposed by combining DL-based RPNs and feature extractors (namely, based on Faster R-CNN or Mask R-CNN) with large-scale Kernel classifiers [22], [23], [24]. According to this approach, the feature extractor is pre-trained off-line on a large representative dataset, yielding a powerful and transferable learned representation, which is kept fixed during on-line operation. The actual regions classification is performed by the integration of an efficient hard-negatives bootstrapping approach (the Minibootstrap [21]) with a set of Kernel-based FALKON classifiers [22], [23].

Labeling efficiency. Labeling efficiency is another key requirement for robotic object detection. The broad class of WSL methods [9], [6] provides a rich set of tools towards this goal in general purpose Computer Vision, in particular

AL and SSL – introduced in Sec. II. After successful applications to deep object classification [25], [26], [27], AL has been recently applied also to object detection [28], [29], [30]. For instance, recently, detection-specific image scoring functions (like e.g., *localization tightness* and *stability* [31]) have been proposed. Instead, when no further annotation is allowed to exploit the unsupervised samples, SSL techniques can be used. Similarly to AL, also SSL has been recently applied to object detection. For instance, in [32], SSL is employed for dataset augmentation and training object detectors. Moreover, the authors point out that vanilla SSL can degrade accuracy in presence of domain shift. We also observed the same issue in our robotic setting and we propose a simple yet effective solution in Sec IV-D. Recent approaches integrate both AL and SSL techniques into the same detection pipeline, such as Self-supervised Sample Mining [33], [34] (SSM). SSM sorts unsupervised images into separate candidate sets for further AL and SSL processing, according to the predictive confidence scores of the underlying DL-based detection model [13]. Another related field in Computer Vision is *Unsupervised Domain adaptation* for object detection [35]. Specifically, the *pseudo-labeling* approach [36], [37] proposes to adapt a detection model to novel and unknown domains (i.e., datasets) by using confident model predictions as pseudo-ground truth.

The aforementioned approaches have been proposed and benchmarked on general purpose Computer Vision datasets. However evaluation of WSL techniques on robotic scenarios is still at an initial stage (e.g., see [7], [5]). For instance, in [5], SSM is extended to enable on-line adaptive object detection for Robotics, by integrating the WSL sample selection strategy with the on-line object detection method [21]. However, [5] still requires a relatively large number of manual annotations, does not investigate the effect of severe domain shift in self-supervision and focuses on a *pool-based* processing. While showing encouraging results, all these limitations prevent its adoption in on-line applications. In this work, we present an empirical analysis of different general purpose computer vision AL and SSL techniques in a challenging robotic scenario, targeting a low annotation budget regime. We focus on how WSL techniques can be used to exploit the robot interaction with the environment and the human teacher to update and improve performance of object detection models previously trained with data of handheld objects. Moreover, we propose solutions to overcome the aforementioned limitations, improving the AL performance and addressing the SSL failure cases under domain shift, increasing overall labeling efficiency.

IV. METHODS

In this work, a robot is asked to detect a set of object instances in an unconstrained environment (referred to as TARGET). A first detection model is trained during a brief interaction with a human, in a teacher-learner scenario, like e.g. in [4] where objects are handheld (the TARGET-LABELED). Then, the robot autonomously explores the environment, acquiring a stream of images in a new setting,

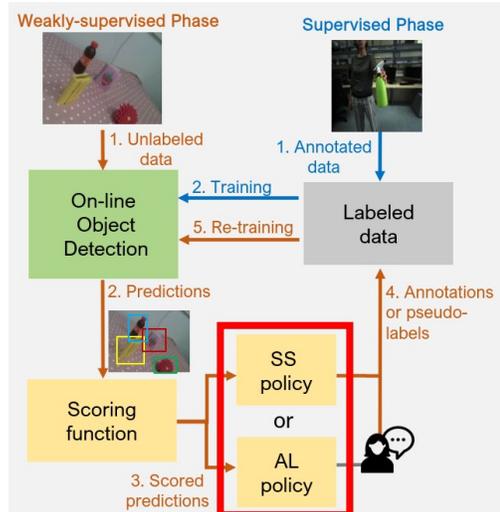


Fig. 1. Overview of the proposed pipeline. Refer to Sec. IV-A for details.

where automatic annotation is not possible. Therefore, these images are not labeled (TARGET-UNLABELED) and are used to adapt the detector on-line exploiting the robot interaction with the environment and the human teacher. In the next sections, we present the proposed pipeline (Sec. IV-A) and the learning protocol (Sec. IV-B). Then, we present all the considered AL and SSL techniques and the proposed approaches (Sec. IV-C and IV-D, respectively).

A. Pipeline Description

The proposed WSL pipeline (see Fig. 1) is composed of four main modules: (i) the *On-line Object Detection*, (ii) the *Scoring function*, (iii) the *AL Selection policy*, and (iv) the *SS Selection policy*.

On-line Object Detection (OOD). For this module, we follow the method proposed in [21], but considering the implementation presented in [38] and [39]. This is an on-line learning approach consisting of two stages: (i) region proposals and feature extraction, and (ii) region classification and bounding-box refinement. The first stage relies on layers from Mask R-CNN [1] (specifically, the convolutional layers, the RPN [40] and the RoI Align layer [1]). In particular, this part is used to extract a set of Regions of Interest (RoIs) from an image and encode them into a set of features. The second stage is composed of a set of FALKON [22] binary classifiers (one for each class of the TARGET) and Regularized Least Squares (RLS) [41], respectively for the classification and the refinement of the proposed RoIs. Classifiers are trained with an approximate bootstrapping approach, called Minibootstrap [21], which addresses the well-known issue of background-foreground class imbalance in object detection [42], while maintaining a short training time. In this work, the adoption of OOD permits to achieve a convenient speed/accuracy trade-off, since it allows to maintain a competitive accuracy with other

DCNN-based approaches with a fraction of the optimization time required (seconds or minutes) [21], [38].

Scoring function. This function assigns a confidence score to the predictions for the images in the TARGET-UNLABELED. This score is then used by the AL and SS Selection policies to decide which images need to be manually annotated or can be considered as pseudo-ground truth. For this part, we employ the Cross-Image Validation (CIV) proposed in SSM [33]. CIV stitches predicted image patches from the TARGET-UNLABELED on random images, sampled from TARGET-LABELED. Then, it executes the detector on the stitched images and computes a *consistency score* from the obtained confidence scores [33].

AL and SS Selection policies. Given the predicted detections obtained by the OOD and the *consistency score* computed by the Scoring function, these two policies decide whether an image of the TARGET-UNLABELED is queried for annotation or the predicted detections are confident enough to be used for self-supervision. Our main contribution relies on these last two components. Firstly, we target a stream-based scenario, since it is more suitable for on-line applications. Secondly, we consider a robotic setting with low annotation budget and a large domain shift of the TARGET-UNLABELED with respect to the TARGET-LABELED. Specifically, for the *AL Selection policy*, we consider several AL techniques, comparing their performance on the considered robotic setting and proposing a solution to enforce diversity during sampling. The adopted AL baselines and the proposed solution are listed in Sec. IV-C. Instead, for the *SS Selection policy*, we consider a stream-based baseline and a novel strategy to overcome issues caused by the domain shift, both described in Sec. IV-D. Finally, another major difference with respect to previous work [5] is that we consider the case in which only one pass over the TARGET-UNLABELED data is allowed, while typically in standard Computer Vision, and also in [5], an iterative process is used. This aspect is crucial for speeding up WSL. However, it makes detector refinement more challenging.

B. Learning Protocol

The learning process is divided into: (i) *Supervised phase* (represented by the light blue arrows in Fig. 1), and (ii) *Weakly-supervised phase* (represented by the orange arrows in Fig. 1). Both phases rely on pre-trained Mask R-CNN’s weights as feature extractor for the OOD. Those weights remain fixed, while model training and adaptation is performed by optimizing on the new data only the second stage of the OOD, i.e., (i) the FALKON classifiers with the Minibootstrap technique and (ii) the RLS box-refinement model (see Sec. IV-A for details). The *Supervised phase* is performed within a few seconds of interaction with a human on the TARGET-LABELED, yielding a first detection model (the *seed model*). In this phase the human shows the objects to the robot, handling them in their hand and annotations are automatically collected. Then, in the *WSL phase*, the SSL

pseudo-ground truth and AL queries are selected from the TARGET-UNLABELED as described in Sec. IV-A, using the *seed model’s* confidence scores. Finally, they are added to the dataset which is used to re-train the on-line detector.

C. Active Learning Strategies

For AL selection, we considered both (i) stream-based approaches, which are the focus of this work, being suited to robotic scenarios, and (ii) pool-based ones.

A simple, yet often effective, pool-based strategy is to sample uniformly at random the images with a confidence score below a threshold (`Uniform random` in Sec. V). Another diversity sampling strategy is to execute *k*-means clustering [41] on the image-level features and select the resulting cluster centers (`K-means-based AL` in Sec. V). In our analysis, we report results for both strategies.

In stream-based AL settings, a simple selection strategy involves confidence score thresholding followed by coin flipping [33] for implementing uncertainty and diversity sampling, respectively (`coin-flip AL` in Sec. V). Another, more practical, solution is to exploit temporal coherence in image sequences to enforce sampling diversity [43]. Leveraging temporal coherence is particularly suitable for on-line robotic tasks, since data, coming in streams, needs to be acquired sequentially and is therefore highly temporally correlated. To this aim, we consider the `Fixed temporal window` strategy, which employs a temporal window of fixed size Δ so that if frame t is selected, any other frame within $[t-\Delta, t+\Delta]$ can no longer be considered for selection. While enforcing diversity, this strategy, by using a fixed Δ , does not take into account: (i) the exploration session duration, that is, the size n_U of TARGET-UNLABELED, which might be known a-priori even in stream-based scenarios, and (ii) the available manual annotation budget k . We show in Sec. V-B that this results in poor performance for low k when the TARGET-UNLABELED is redundant.

To overcome this limitation, we propose to use an adaptive temporal window size, defined as

$$\Delta_{n_U, k} = \frac{n_U \cdot \alpha}{k}$$

and referred to as `Adaptive temporal window` in Sec. V. This strategy allows to tailor the strictness (window size) of the temporal diversity-enforcing sampling to the overall amount of available unsupervised data n_U , while at the same time ensuring to make full use of the available budget k . For instance, given a budget k , the adaptive window size grows linearly with n_U in order to cover the entire duration of the exploration session. $\alpha \in (0, k) \subseteq \mathbb{R}$ is a hyperparameter accounting for the proportion of AL candidates with respect to n_U , which is unknown a priori.

D. Semi-supervised Learning Strategies

For SS selection, we consider two stream-based baselines. The first is the `SS baseline`, which selects all the images passing CIV as pseudo ground truth. However, we show in Sec. V-C that under domain shift this leads to model degradation due to the abundance of false negatives. For

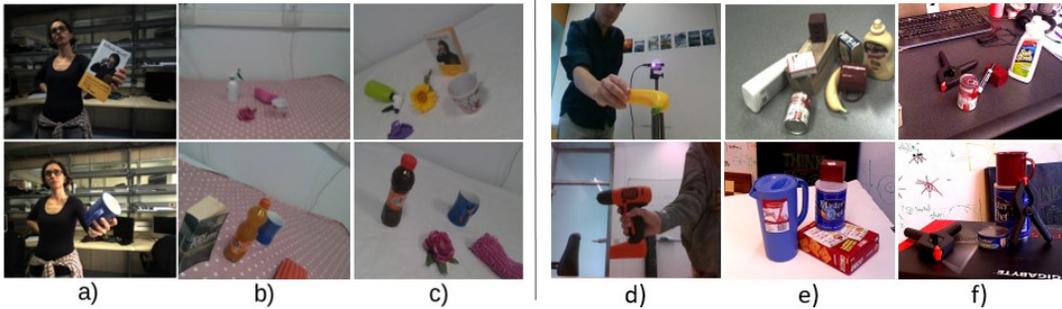


Fig. 2. Example images of the datasets used for this work: **a)** iCWT dataset; **b)** POIS in TABLE-TOP dataset; **c)** WHITE in TABLE-TOP dataset; **d)** HO-3D; **e)** YCB-Video training set, **f)** YCB-Video test set.

this reason, in this work, we propose a more conservative strategy, namely $SS_{pos. only}$, which only selects positive predictions and leaves out negative ones. In Sec. V-C, we show that our approach successfully counteracts severe model degradation.

V. EXPERIMENTS

The objective of our experiments is to evaluate the performance of the presented WSL techniques in improving detection performance under domain shift. Specifically, we consider the scenario of a robot previously trained with human interaction to detect handheld objects. We aim to generalize to a different setting (i.e., a table top) by exploiting the unlabeled data collected by the robot during autonomous exploration (*Weakly-supervised phase* in Sec. IV-B).

A. Experimental Setup

For the OOD, the weights of the Feature extractor are learned by training Mask R-CNN on the MS COCO [44] dataset. ResNet50 [45] has been considered as Mask R-CNN’s convolutional backbone (we use the available pre-trained Mask R-CNN weights¹). During *Supervised* and *Weakly-supervised phases*, the feature extractor is fixed, while the FALKON classifiers and RLS are updated as explained in Sec. IV-B (we relied on [21] for hyperparameters selection). This allows to achieve a training time of few seconds or minutes for each learning step.

Given the aforementioned target scenario, in our experiments we consider two different cases of domain adaptation from handheld (*Supervised phase*) to table-top objects (*Weakly-supervised phase*). Specifically, we adapt (i) from iCubWorld Transformations [3] (iCWT) to a set of sequences depicting a subset of iCWT’s objects on a table-top (TABLE-TOP) and (ii) from HO-3D [46] to YCB-Video [47].

From iCWT to TABLE-TOP (iCubWorld domain). iCWT contains images for 200 handheld objects. Each object is demonstrated by a human teacher to the robot (as in [4]) and is acquired with different sequences representing specific viewpoint transformations: 2D rotation (2D ROT), generic

rotation (3D ROT), translation (TRANSL), scaling (SCALE) and all transformations (MIX) (see [3]). For the *Supervised phase*, we employ a subset of the iCWT, considering 21 of the total 200 objects. All the transformations, except from MIX, are considered, resulting in a TARGET-LABELED of size $n_L \sim 6K$. The TABLE-TOP depicts the same 21 objects randomly placed on a table with two different tablecloths: (i) pink/white pois (POIS) and (ii) white (WHITE). The two datasets contain the same objects, but with an important domain shift: iCWT frames include the hand of the teacher, whereas TABLE-TOP has different backgrounds and light conditions and depicts objects on a table. Refer to Fig. 2 for a visual representation of the domain shift. For the *Weakly-supervised phase*, we consider the WHITE sequence as TARGET-UNLABELED while we leave the POIS sequence as test set to evaluate performance. These two sets are respectively of size $\sim 2K$ and $\sim 1K$.

From HO-3D to YCB-Video (YCB domain). Similarly, in HO-3D and YCB-Video, objects from the YCB [48] dataset are presented handheld by a human and in table-top sequences, respectively. Specifically, YCB-Video presents sequences for 21 objects while in HO-3D a subset of 9 of those objects are considered. Note that for our experiments we do not consider the labels for the remaining 12 objects in YCB-Video. For the *Supervised phase*, we take from HO-3D at most four sequences for each object resulting in a TARGET-LABELED of size $n_L \sim 20K$. For the *Weakly-supervised phase*, we consider a set of $\sim 11.3K$ frames obtained by extracting one image every ten from the total 80 training video sequences available in the YCB-Video. As test set, instead, we consider the $\sim 3K$ keyframe [47] images chosen from the remaining 12 sequences in the YCB-Video.

Evaluation metrics. We report performance in terms of mAP (mean Average Precision) at the IoU (Intersection over Union) threshold set to 0.5, as defined for Pascal VOC 2007 (see [49]). Specifically, we repeat each experiment for three trials and we present the results, reporting the mean and the standard deviation of the obtained accuracy².

¹https://github.com/facebookresearch/maskrcnn-benchmark/blob/master/MODEL_ZOO.md

²All experiments have been executed on a machine equipped with Intel Xeon E5-2690 v4 CPUs @2.60GHz, and an NVIDIA Tesla P100 GPU.

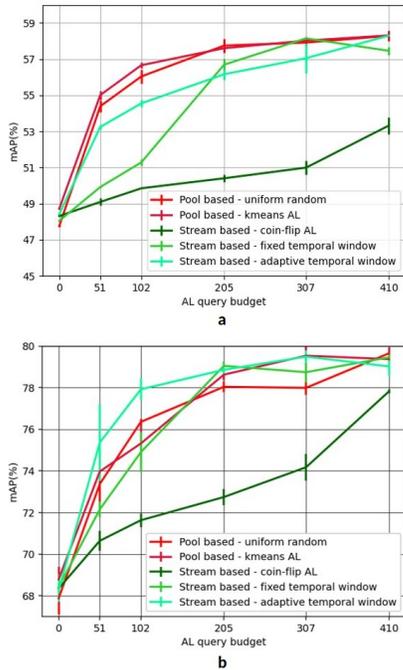


Fig. 3. mAP comparison of pool-based (red) and stream-based (green) AL strategies with varying query budgets for iCubWorld (a) and YCB (b) domains.

B. Active Learning Sampling Strategy Evaluation

In this section, we compare the AL techniques described in Sec. IV considering different manual annotation budgets for both iCubWorld and YCB domains. To this aim, we report in Fig. 3a and b the mAP trends obtained by increasing the AL query budget during the *Weakly-supervised phase*. Specifically, we report in red shades the performance obtained by the pool-based strategies (namely, *k*-means-based AL and Uniform random from Sec. IV), and in green shades the stream-based ones (namely, Coin-flip AL, Fixed temporal window, and the proposed Adaptive temporal window from Sec. IV). We empirically set the fixed temporal window size as $\Delta = 6$ and the adaptive temporal window hyperparameter as $\alpha = 0.5$ for the iCubWorld domain and $\alpha = 0.4$ for the YCB domain. As it can be observed in Fig. 3a, for iCubWorld domain the pool-based methods achieve the best mAP trends. Notably, we observe that the Uniform random baseline is almost as effective as *k*-means based-AL and they both present an early steep slope for limited manual annotation budgets. These two aspects are due to the fact that the considered TABLE-TOP dataset in the iCubWorld domain, contains sequences of similar (and thus redundant) frames which need to be properly filtered during data selection. This aspect of the dataset is also the main cause for the poor performance obtained by the two stream-based techniques: Coin-flip AL and Fixed temporal window, for low numbers of manual annotations. Indeed, while being more suited for a robotic application, by reasoning only on a frame-by-frame fashion, they lack global information on

TABLE I
RESULTS OBTAINED BY SS BASELINE (2nd COLUMN) AND SS POSITIVES (3rd COLUMN) FOR LARGE (1st ROW) AND SMALL (2nd ROW) DOMAIN SHIFT FROM THE SUPERVISED (1st COLUMN) TO THE WEAKLY-SUPERVISED PHASE.

	Sup. phase (mAP(%))	SS baseline (mAP(%))	SS pos. only (mAP(%))	SS samples
Large DS	48.8 ± 0.3	37.9 ± 1.8	50.9 ± 0.06	$\sim 12\%$
Small DS	40.3 ± 0.9	47.1 ± 0.1	46.6 ± 0.2	$\sim 35\%$

the whole data distribution, which turns out to be a critical drawback especially for limited manual annotation budgets. However, for higher budgets, the Fixed temporal window baseline achieves accuracies closer to the pool-based ones. Finally, the proposed Adaptive temporal window presents the best mAP trend, among the stream-based approaches, especially for low annotation budgets and it closely matches the pool-based ones. On the contrary, as it can be observed in Fig. 3b, for the YCB domain the pool-based methods, the Fixed and the Adaptive temporal window present similar mAP trends. Specifically, the proposed Adaptive temporal window has the steepest slope. This is due to the fact that the YCB-Video dataset, differently from the TABLE-TOP, presents less redundant sequences and a smarter data selection based on temporal coherence provides the best performance.

Finally, it is important to note that the proposed Adaptive temporal window stream-based approach achieves significantly higher mAP values, than other stream-based techniques, for low annotation budgets for both domains. This makes it the most successful stream-based approach in such regime, which is the target of the presented work.

C. Semi-supervised Learning Evaluation

In this section, we investigate the impact of domain shift from handheld objects to table-top datasets when no labeling is allowed (i.e., SSL). To this aim, we report in Tab. I the results of applying the SS baseline (as defined in Sec. IV) in the two following scenarios:

- **Large domain shift.** In this case (*Large DS* row in Tab. I), we consider the scenario in which the TARGET-UNLABELED presents a completely different setting (i.e., a table top) with respect to the TARGET-LABELED (i.e., hand-held). To this end we used the two datasets described in Sec. V-A.
- **Small domain shift.** In this case (*Small DS* row in Tab. I), TARGET-LABELED and TARGET-UNLABELED present similar conditions. The only difference in the latter one is that the objects are presented, unlabeled, with different view poses. To this aim, we considered as TARGET a 30-object identification task from iCWT. For each object, we then use the TRANSL sequence ($\sim 2K$ images) as TARGET-LABELED and the union of the 2D ROT, 3D ROT, and SCALE sequences ($\sim 6K$ images) as the TARGET-UNLABELED. We test on the MIX sequences of all the objects ($\sim 4.5K$ images).

Note that, in this experiment, we use the iCubWorld domain only because the explicit sub-division in different viewpoint transformations of iCWT allows to control the dataset split in TARGET-LABELED and TARGET-UNLABELED such that they present similar, but not identical, conditions. This allows to precisely identify the *Small DS* setting.

Tab. I reports the results obtained in both cases. For each row, we report the mAP (represented as mean and standard deviation of the different repetitions) after the *Supervised phase* (first column) and after the *Weakly-supervised phase* for both *SS baseline* and *SS pos. only* (second and third columns). Moreover, in the fourth column we report the average percentage of samples selected by the SS process over the total. As it can be observed, adding self-supervised data with *SS baseline*, with small domain shift, results in an improvement in accuracy. On the contrary, with a larger domain shift, it leads to a significant accuracy deterioration. A reason for this phenomenon can be identified by analyzing the pseudo-ground truth generated by the SS process. We report in Fig. 4 some representative images depicting in green the region proposal candidates classified as background by the detection system and that are therefore added as negative samples to the dataset by the *SS baseline*. The actual detections which instead are considered as positive samples in the SS process are shown in red. It can be noticed that, with large domain shift, only few objects are correctly detected and therefore added to the training set as positives, while most others are false negatives which are automatically annotated as background samples. Clearly, retraining the detection model with such a poorly-labeled dataset leads to the sharp performance decay shown in Tab. I. This confirms similar findings from the literature [36], in the considered setting. Note that, we empirically noticed that lowering the confidence threshold used to determine a positive prediction is not suitable since, while not ensuring less false negatives, it leads to imprecise predictions, with a similar negative effect on the subsequent training. In Sec. IV, we introduce the *SS pos. only* to address this issue. This more conservative strategy includes only the regions predicted as positive in the SS dataset, while the others are filtered away, avoiding adding false negatives. Third column in Tab. I shows that this strategy, does not modify the baseline in case of Small DS where SS data is already reliable. However, for Large DS, not only effectively removes wrong labels from the dataset, recovering from the two-digits accuracy decay, but also successfully yields a performance improvement of ~ 2 points. Moreover, it allows to drastically reduce the standard deviation of the obtained accuracy, from 1.8% to 0.06%, being less sensitive to statistical fluctuations. This demonstrates that, in cases when no human manual annotation is allowed, a robot trained to detect handheld objects can explore the new domain, self-annotating the newly collected data and improving detection performance.

VI. CONCLUSIONS

In this paper, we target the scenario of a robot trained with human interaction to detect handheld objects, aiming

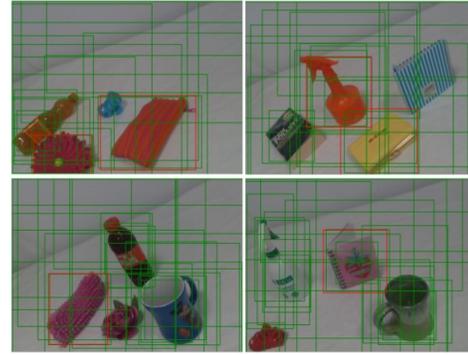


Fig. 4. Example predictions on the TARGET-UNLABELED before model adaptation, selected by the *SS Baseline* for SS as positives (red boxes) and negatives (green boxes).

to improve detection performance in different settings with autonomous exploration and limited human intervention. We empirically demonstrate that general purpose WSL techniques are unsuitable for challenging robotic scenarios and we propose solutions to both (i) enforce diversity sampling for AL queries and (ii) improve strong positives selection for SSL under severe domain shift. Finally, we build on previous work [5], presenting and empirically evaluating a stream-based weakly-supervised on-line object detection pipeline for Robotics, which exploit the robot interaction with the environment and the human teacher to update and improve performance of the visual system. It significantly alleviates the annotation burden for on-line model adaptation to novel settings while maximizing accuracy.

REFERENCES

- [1] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, “Mask r-cnn,” *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 2980–2988, 2017.
- [2] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, “Yolov4: Optimal speed and accuracy of object detection,” *arXiv preprint arXiv:2004.10934*, 2020.
- [3] G. Pasquale, C. Ciliberto, F. Odone, L. Rosasco, and L. Natale, “Are we done with object recognition? the icub robot’s perspective,” *Robotics and Autonomous Systems*, vol. 112, pp. 260 – 281, 2019. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0921889018300332>
- [4] E. Maiettini, G. Pasquale, L. Rosasco, and L. Natale, “Interactive data collection for deep learning object detectors on humanoid robots,” in *2017 IEEE-RAS 17th International Conference on Humanoid Robotics (Humanoids)*, Nov 2017, pp. 862–868.
- [5] E. Maiettini, G. Pasquale, V. Tikhonoff, L. Rosasco, and L. Natale, “A weakly supervised strategy for learning object detection on a humanoid robot,” in *2019 IEEE-RAS 19th International Conference on Humanoid Robotics (Humanoids)*, Nov 2019.
- [6] Z.-H. Zhou, “A brief introduction to weakly supervised learning,” *National Science Review*, vol. 5, no. 1, pp. 44–53, Jan. 2018, publisher: Oxford Academic. [Online]. Available: <https://academic.oup.com/nsr/article/5/1/44/4093912>
- [7] S. Jamieson, J. P. How, and Y. Girdhar, “Active reward learning for robotic vision based exploration in bandwidth limited environments,” in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, 2020, pp. 1806–1812.
- [8] B. Settles, *Active Learning*, ser. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool, 2012.
- [9] J. Hernández-González, I. Inza, and J. A. Lozano, “Weak supervision and other non-standard classification problems: A taxonomy,” *Pattern*

- Recognition Letters*, vol. 69, pp. 49–55, Jan. 2016. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0167865515003505>
- [10] C. Papageorgiou, M. Oren, and T. Poggio, “A general framework for object detection,” in *Sixth International Conference on Computer Vision (IEEE Cat. No.98CH36271)*, Jan. 1998, pp. 555–562.
 - [11] P. Viola and M. Jones, *Rapid object detection using a boosted cascade of simple features*, 2001.
 - [12] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks,” in *Advances in Neural Information Processing Systems 28*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds. Curran Associates, Inc., 2015, pp. 91–99.
 - [13] J. Dai, Y. Li, K. He, and J. Sun, “R-FCN: Object Detection via Region-based Fully Convolutional Networks,” in *Advances in Neural Information Processing Systems 29*, D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, Eds. Curran Associates, Inc., 2016, pp. 379–387.
 - [14] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, and Y. Wei, “Deformable convolutional networks,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 764–773.
 - [15] J. Pang, K. Chen, J. Shi, H. Feng, W. Ouyang, and D. Lin, “Libra R-CNN: Towards Balanced Learning for Object Detection,” in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019, pp. 821–830, iSSN: 2575-7075.
 - [16] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, “SSD: Single Shot MultiBox Detector,” in *Computer Vision – ECCV 2016*, ser. Lecture Notes in Computer Science, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds. Cham: Springer International Publishing, 2016, pp. 21–37.
 - [17] S. Zhai, D. Shang, S. Wang, and S. Dong, “Df-ssd: An improved ssd object detection algorithm based on densenet and feature fusion,” *IEEE Access*, vol. 8, pp. 24 344–24 357, 2020.
 - [18] T. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, “Focal Loss for Dense Object Detection,” in *2017 IEEE International Conference on Computer Vision (ICCV)*, Oct. 2017, pp. 2999–3007, iSSN: 2380-7504.
 - [19] I. J. Goodfellow, M. Mirza, D. Xiao, A. Courville, and Y. Bengio, “An empirical investigation of catastrophic forgetting in gradient-based neural networks,” *arXiv preprint arXiv:1312.6211*, 2013.
 - [20] E. Maittini, G. Pasquale, L. Rosasco, and L. Natale, “Speeding-up object detection training for robotics with falkon,” in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct 2018.
 - [21] —, “On-line object detection: a robotics challenge,” *Autonomous Robots*, Nov 2019. [Online]. Available: <https://doi.org/10.1007/s10514-019-09894-9>
 - [22] A. Rudi, L. Carratino, and L. Rosasco, “Falkon: An optimal large scale kernel method,” in *Advances in Neural Information Processing Systems*, 2017, pp. 3888–3898.
 - [23] G. Meanti, L. Carratino, L. Rosasco, and A. Rudi, “Kernel methods through the roof: handling billions of points efficiently,” 2020, eprint: 2006.10350.
 - [24] A. Rudi, R. Camoriano, and L. Rosasco, “Less is more: Nyström computational regularization,” in *Advances in Neural Information Processing Systems*, 2015, pp. 1657–1665.
 - [25] A. Kirsch, J. van Amersfoort, and Y. Gal, “BatchBALD: Efficient and Diverse Batch Acquisition for Deep Bayesian Active Learning,” *arXiv:1906.08158 [cs, stat]*, Oct. 2019, arXiv: 1906.08158. [Online]. Available: <http://arxiv.org/abs/1906.08158>
 - [26] F. Zhdanov, “Diverse mini-batch Active Learning,” *arXiv:1901.05954 [cs, stat]*, Jan. 2019, arXiv: 1901.05954. [Online]. Available: <http://arxiv.org/abs/1901.05954>
 - [27] J. T. Ash, C. Zhang, A. Krishnamurthy, J. Langford, and A. Agarwal, “Deep Batch Active Learning by Diverse, Uncertain Gradient Lower Bounds,” Jan. 2020. [Online]. Available: <https://openreview.net/forum?id=0HjEAtQNNWD>
 - [28] H. H. Aghdam, A. Gonzalez-Garcia, A. Lopez, and J. Weijer, “Active Learning for Deep Detection Neural Networks,” in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. Seoul, Korea (South): IEEE, Oct. 2019, pp. 3671–3679. [Online]. Available: <https://ieeexplore.ieee.org/document/9009535/>
 - [29] E. Haussmann, M. Fenzi, K. Chitta, J. Ivanecky, H. Xu, D. Roy, A. Mittel, N. Koumchatzky, C. Farabet, and J. M. Alvarez, “Scalable Active Learning for Object Detection,” *arXiv:2004.04699 [cs]*, Apr. 2020, arXiv: 2004.04699. [Online]. Available: <http://arxiv.org/abs/2004.04699>
 - [30] S. V. Desai, A. C. Lagandula, W. Guo, S. Ninomiya, and V. N. Balasubramanian, “An Adaptive Supervision Framework for Active Learning in Object Detection,” 2019.
 - [31] C.-C. Kao, T.-Y. Lee, P. Sen, and M.-Y. Liu, “Localization-Aware Active Learning for Object Detection,” in *Computer Vision – ACCV 2018*, ser. Lecture Notes in Computer Science, C. Jawahar, H. Li, G. Mori, and K. Schindler, Eds. Cham: Springer International Publishing, 2019, pp. 506–522.
 - [32] Y. Li, D. Huang, D. Qin, L. Wang, and B. Gong, “Improving Object Detection with Selective Self-supervised Self-training,” in *European Conference on Computer Vision*. Springer, 2020, pp. 589–607.
 - [33] K. Wang, X. Yan, D. Zhang, L. Zhang, and L. Lin, “Towards Human-Machine Cooperation: Self-supervised Sample Mining for Object Detection,” *arXiv:1803.09867 [cs]*, Mar. 2018, arXiv: 1803.09867. [Online]. Available: <http://arxiv.org/abs/1803.09867>
 - [34] K. Wang, L. Lin, X. Yan, Z. Chen, D. Zhang, and L. Zhang, “Cost-effective object detection: Active sample mining with switchable selection criteria,” *IEEE transactions on neural networks and learning systems*, vol. 30, no. 3, pp. 834–850, 2019, publisher: IEEE.
 - [35] P. Oza, V. A. Sindagi, V. VS, and V. M. Patel, “Unsupervised domain adaptation of object detectors: A survey,” *arXiv preprint arXiv:2105.13502*, 2021.
 - [36] X. Li, W. Chen, D. Xie, S. Yang, P. Yuan, S. Pu, and Y. Zhuang, “A free lunch for unsupervised domain adaptive object detection without source data,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 10, 2021, pp. 8474–8481.
 - [37] A. RoyChowdhury, P. Chakrabarty, A. Singh, S. Jin, H. Jiang, L. Cao, and E. Learned-Miller, “Automatic adaptation of object detectors to new domains using self-training,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
 - [38] F. Ceola, E. Maittini, G. Pasquale, L. Rosasco, and L. Natale, “Fast region proposal learning for object detection for robotics,” *arXiv preprint arXiv:2011.12790*, 2021.
 - [39] —, “Fast object segmentation learning with kernel-based methods for robotics,” in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, 2021.
 - [40] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards real-time object detection with region proposal networks,” in *Neural Information Processing Systems (NIPS)*, 2015.
 - [41] T. Hastie, R. Tibshirani, and J. Friedman, *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2009.
 - [42] K. Oksuz, B. C. Cam, S. Kalkan, and E. Akbas, “Imbalance problems in object detection: A review,” *IEEE transactions on pattern analysis and machine intelligence*, 2020.
 - [43] H. H. Aghdam, A. Gonzalez-Garcia, J. v. d. Weijer, and A. M. Lopez, “Active learning for deep detection neural networks,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
 - [44] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *European Conference on Computer Vision (ECCV)*, Zürich, 2014, oral. [Online]. Available: <http://arxiv.org/abs/1411.1724>
 - [45] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” *arXiv preprint arXiv:1512.03385*, 2015.
 - [46] S. Hampali, M. Rad, M. Oberweger, and V. Lepetit, “HONotate: A method for 3D annotation of hand and object poses,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3196–3206.
 - [47] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, “PoseCNN: A convolutional neural network for 6d object pose estimation in cluttered scenes,” 2018.
 - [48] B. Calli, A. Singh, A. Walsman, S. Srinivasa, P. Abbeel, and A. M. Dollar, “The YCB object and model set: Towards common benchmarks for manipulation research,” in *2015 international conference on advanced robotics (ICAR)*. IEEE, 2015, pp. 510–517.
 - [49] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, “The Pascal Visual Object Classes (VOC) Challenge,” *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, June 2010. [Online]. Available: <https://doi.org/10.1007/s11263-009-0275-4>