A Novel Dataset for Evaluating and Alleviating Domain Shift for Human Detection in Agricultural Fields

Paraskevi Nousi¹, Emmanouil Mpampis¹, Nikolaos Passalis¹, Ole Green², Anastasios Tefas¹ ¹Computational Intelligence and Deep Learning Group, Department of Informatics,

Aristotle University of Thessaloniki, Thessaloniki, Greece

²Agrointelli, Aarhus, Denmark

E-mails: {paranous,empampis,passalis}@csd.auth.gr, olg@agrointelli.com, tefas@csd.auth.gr

Abstract-In this paper we evaluate the impact of domain shift on human detection models trained on well known object detection datasets when deployed on data outside the distribution of the training set, as well as propose methods to alleviate such phenomena based on the available annotations from the target domain. Specifically, we introduce the OpenDR Humans in Field dataset, collected in the context of agricultural robotics applications, using the Robotti platform, allowing for quantitatively measuring the impact of domain shift in such applications. Furthermore, we examine the importance of manual annotation by evaluating three distinct scenarios concerning the training data: a) only negative samples, i.e., no depicted humans, b) only positive samples, i.e., only images which contain humans, and c) both negative and positive samples. Our results indicate that good performance can be achieved even when using only negative samples, if additional consideration is given to the training process. We also find that positive samples increase performance especially in terms of better localization. The dataset is publicly available for download at https://github.com/opendr-eu/datasets.

I. INTRODUCTION

Object detection combines the tasks of classification and localization, i.e., it refers to finding *what* objects are pictured in an image, as well as *where* in the image they are located. In the case of multiple object, multiple class object detection, a generic object detector should be able to detect an unknown number of objects belonging to a number of different classes. Depending on the training dataset, these classes can include people, animals, inanimate objects, etc. Such datasets include the widely popular PASCAL VOC [1] and MS COCO [2] object detection benchmarks, containing objects from 20 and 80 classes respectively. Deep Learning brought significant improvements both in terms of effectiveness and efficiency and currently the top performing object detection methods on these challenging benchmarks are all based on Deep Convolutional Neural Networks (CNNs).

Despite improvements in these detectors on well known object detection benchmarks, deploying them on new applications highlights the domain adaptation problem. Domain adaptation refers to the process of alleviating the domain shift between a source and target domain, i.e., how to effectively deploy a detector trained on a source domain onto a target



(a) Image depicting two humans, captured by the front camera of the robot.



(b) Image depicting one human, captured by the back camera of the robot.

Fig. 1: Examples of detections using a pretrained model. Existing models lead to false detections under distribution shifts.

domain, which differs from the source in some way. Recent works in domain adaptation for object detection propose a progressive shift towards the target domain [3], [4]. A somewhat more straightforward approach to domain adaptation is incremental learning [5]. In any case, knowledge transfer is of significant importance when training a detector on a new dataset, in order to maximize its accuracy for the new domain.

In this paper, we introduce a dataset, called OpenDR Humans in Field, collected in the context of agricultural use cases using the Robotti robotic platform, designed for detection of humans in fields. Samples from this dataset are shown in Figure 1, with detections made using an SSD detector [6] pretrained on the COCO dataset [2], where the domain shift problem is evident. We evaluate the accuracy of various detectors trained on existing datasets, identifying important limitations that these detectors face on scenarios like this. The results of this evaluation highlight the need to use domain adaptation and knowledge transfer approaches to increase the performance of detection on such applications. We examine various methods in depth and report results in terms of precision and recall for each. Our main findings are that using only negative samples can significantly drop the false positive rate, compared to a baseline pretrained model, and incorporating positive samples can further improve localization, leading to increased detection precision.

II. RELATED WORK

Single-stage detectors have been shown to perform about as well as their two-stage, heavyweight counterparts, while running at much faster speeds. The seminal methods of YOLO [7] and SSD [6] inspired many recent works which utilize the anchor-based, single-stage architectures proposed by them. Anchor-free object detectors aim to tackle issues arising from the use of predefined anchors, such as the need for thousands of such anchors in order to train dense object detectors, or the tedious hyperparameters they introduce, like the size, aspect ratio etc. CenterNet [8] is one such anchor-free object detector, taking into consideration the center of objects as well as the corners, to detect each object as a triplet.

In the context of agriculture, object detection methods can assist robots in their tasks in various ways [9]. We are specifically interested in the human-centric scenario, where human labour is complemented with robots, focusing on human detection in fields, which is a critical safety aspect for human-robot interaction. The main contribution of this paper is the collection of a dataset that depicts humans in agricultural fields in various conditions. This is in contrast to the most commonly used person detection datasets, where humans are depicted in urban scenarios. Furthermore, the lenses attached to the robot are wide-angled, leading to bounding boxes of different proportions than those commonly seen in existing datasets. Therefore, the collected dataset allows for evaluating the impact of domain shift, as well as employing method for reducing its effect.

Indeed, this domain shift problem [10], compared to existing datasets, is evident in Figure 2, and is encountered in other computer vision tasks as well, such semantic segmentation [11] or concept detection [12]. In object detection, a two-level domain adaptation approach was introduced in [13] for

Faster R-CNN, on an image-level as well as on an instancelevel. In [14], a multidomain-invariant representation learning process was proposed, using adversarial learning. In this work, we tackle the domain shift problem in a data-driven manner, influenced also by the lack of a large collection of images, paving the way for developing methods that can work under the challenging settings that are often encountered in agricultural applications.

The rest of this paper is structured as follows. Section II presents several works related to object detection and domain adaptation. The dataset collection process, as well as the employed methods for alleviating domain shift are described in Section III. The results of our experimental study are presented and analyzed in Section IV. Finally, Section V concludes our work and summarizes our findings.

III. PROPOSED METHOD

A. Dataset Collection

A Robotti was deployed by AGI to collect images with a front and back camera, in a realistic scenario to mimic the images that the robot might encounter in the agricultural use case. A total of 8038 images were collected on two separate occasions, 818 in the first batch and 7233 in the second. For the purposes of this work, the first batch was fully annotated, while the second one is provided to support unsupervised learning tasks. Of the 818 collected images, 13 were discarded as they depicted unwilling participants to comply with GDPR. The remaining images were annotated with bounding boxes, where one bounding box corresponds to one depicted person. The LabelImg¹ tool was used for the annotation, which outputs annotations in PASCAL VOC .xml format. Figure 2a is an image from this dataset annotated with two bounding boxes for the two depicted humans. In total, 158 images contained people, and 647 images did not. The latter were annotated with an empty bounding box list, to be used as negative samples in object detection algorithms. Figure 2b is an example of an image from this dataset containing no humans.

B. Alleviating the Domain Shift

A natural first step to alleviate the domain shift is to finetune a pretrained detector on background images from the new domain, i.e., images which do not depict any objects of interest. One benefit of this method is that no annotation is required, which can often be a tedious and time-consuming activity. However, training a detector solely on negative samples may quickly degrade the detector's performance on positive samples, i.e., images depicting objects of interest. A small learning rate and only a small number of training iterations can be used to lessen this undesirable side-effect.

In the absence of a sufficient number of negative training examples, or in the case where performance is still subpar, positive samples should be included in the training set. Although bounding box annotation is costly, it is the most reliable way to increase detection performance on a new

¹https://github.com/tzutalin/labelImg



(a) Positive sample depicting two humans.



(b) Negative sample with no humans in field of view.

Fig. 2: Examples of collected images: (a) Image depicting two humans, annotated with bounding boxes, (b) Image depicting no humans.

dataset. Positive samples can be used as the only training set as they contain some of the background information as well. Furthermore, finetuning a pretrained detector with both positive and negative samples can lead to fewer false positive detections in comparison to training with only positive samples, as well as fewer missed detections in comparison to using only negative samples.

However, finetuning using only the target dataset can deteriorate the detector's performance on the source dataset due to catastrophic forgetting phenomena [15]. Even though the source dataset may no longer be relevant, this drop in performance can be reflected in the target dataset, in the form of overfitting. Thus, we propose that in addition to the aforementioned training sets, the source dataset is used in finetuning as well. Specifically for human detection, only the 'person' class from the source dataset is extracted and appended to the training set. To enforce balance between the samples of the source and target datasets, the samples of the latter are repeated multiple times. Each sample undergoes transformations, according to the data augmentation protocol of the detector, such that even if the same image is used twice in a batch, the detector sees a slightly different version of it.

TABLE I: Evaluation in terms of precision at 0.5 IoU, of pretrained detectors on the collected dataset depicting humans in field.

Method	Train Set	Pos. Only	All	FPS
SSD	VOC	53.5	42.2	23.6
SSD	COCO	80.3	70.1	23.6
SSD - MBNet	VOC	40.8	18.8	35
SSD - MBNet	COCO	60.7	42.3	35
CenterNet	VOC	43.4	28.6	16.1
CenterNet	COCO	63.1	54.8	16.1
YOLOv3	VOC	63.4	60.9	15.2
YOLOv3	COCO	78.9	74.7	15.2

IV. EXPERIMENTAL RESULTS

A. Baseline models

An extensive evaluation of pretrained detectors of the SSD [6], YOLOv3 [7] and CenterNet [8] families, was conducted for this dataset. The results are summarized in Table I in terms of precision at 0.5 IoU and FPS on Jetson AGX. The detectors are trained on either the PASCAL VOC [1] and MS COCO [2] object detection benchmarks, containing objects of 20 and 80 classes respectively. Finally, the MobileNet version of SSD [16], [17] is also evaluated. For the target dataset, we evaluate the methods on two subsets: a) on the positive samples only ('Pos. Only'), and b) on the entire test set ('All'), including both positive and negative samples. The reason behind this choice is to examine the effect of each training method on the false positive detections.

As expected, the addition of images without people highlights the false positive accumulation, due to the unseen backgrounds present in the dataset. Detectors trained on COCO seem to perform significantly better than those trained on VOC, which can be attributed to the wider range of appearance in people in the larger COCO dataset. The object scale in COCO is also more varied, containing people as small as 10 pixels in height. The YOLOv3 detector in general performs the best, but is the slowest of the evaluated detectors on the Jetson AGX. The SSD MobileNet variant, especially when trained on COCO, seems to give off the best speed/accuracy trade-off. Even so, the drop in precision is significant when considering negative-only samples.

Based on this experimental study, we conclude that further training of the detectors is necessary to avoid false positive detections as well as to increase the true positive ratio. Knowledge transfer from the COCO dataset seems to be the most promising direction, as it leads to the best precision for all detectors. Furthermore, we choose the SSD algorithm as it is the fastest of the compared ones, and specifically the standard VGG16 version, as it still runs at about realtime on the AGX while achieving higher performance than its MobileNet counterpart.

B. Domain Adaptation Experiments

Two major sets of experiments are conducted. In the first case, the detector is finetuned using only the target dataset, and specifically different splits of it. In the second case,

TABLE II: Target domain finetuning - Evaluation using only the positive samples of the dataset

Train Set	AP	Precicion@0.5	Recall@0.5
Baseline (COCO)	49.8	80.3	55.4
Finetune with negatives	50.6	83.1	55.5
Finetune with positives	57.0	90.1	63.2
Finetune with both	56.6	91.1	63.3

TABLE III: Target domain finetuning - Evaluation using the entire (positive and background) samples of the dataset

Train Set	AP	Precicion@0.5	Recall@0.5
Baseline (COCO)	44.8	70.1	55.4
Finetune with negatives	48.6	78.9	55.5
Finetune with positives	56.7	89.3	63.2
Finetune with both	56.5	90.9	63.3

the detector is finetuned using the target dataset as well as the COCO 'person' subset, i.e., any images from the COCO dataset which depict humans.

1) Finetuning on target domain only: For the following experiments, we also measure the performance in terms of Average Precision (AP), Precision at 0.5 IoU threshold, and Recall at 0.5 IoU. Thus we can draw conclusions regarding the false positive (FP), false negative (FN) and localization performance of each method. Table II contains the results of our study on the positive subset of the target dataset. In comparison to the pretrained model on COCO, using only negative samples increases all metrics, although it has the least significant effect on recall. This can be attributed to a very small change of the FN detections, i.e., the detector is only slightly better at finding humans it didn't before finetuning. The most significant change is in terms of Precision@0.5, translating to a smaller FP rate, i.e., the detector has learned to not make false predictions, as expected. In terms of AP, the increase is not as large, indicating that although the overall FP rate has improved, localization issues ensue. This highlights the need to add positive samples to the training set.

On the other hand, using only positive samples, significantly increases the detection performance in terms of both precision and recall. Using both positive and negative samples further increases the precision at 0.5 IoU, at the cost of slightly worsened localization at higher thresholds. Furthermore, the effect on FN is negligible. This indicates that using positive only samples provides the detector with enough background (i.e., negative) samples to reach this peak performance.

The performance of the proposed methods on the entire test set (both positive and negative samples) is shown in Table III. Note that splitting the test set like this only affects the precision scores, and not the recall. Thus, we focus on the precision scores, and specifically on the FP and localization performance.

All of the evaluated methods lead to more FP detections, which is expected as these occur on the added negatives-only subset. Other than this drop, the results are similar to those

TABLE IV: Finetuning using both the target and source domain - Evaluation using the entire (positive and background) samples of the dataset

Train Set	AP	Prec @0.5	Rec @0.5	COCO AP
Baseline (COCO)	44.8	70.1	55.4	37.0
COCO+Neg.	49.1	79.4	53.5	
COCO+Pos.	61.9	94.3	67.0	34.9
COCO+Both	60.7	93.9	70.1	36.8

regarding the positives-only subset. Specifically, using negative samples only improves the baseline performance and the effect is more prominent on this set (+8.8% precision at 0.5 IoU, in comparison to +2.8% in positives-only).

Using only positive and using both positive and negative samples both significantly improve the performance over the baseline pretrained model, and actually more or less reach the same performance as when evaluating only on positive samples. This result is consistent with the fact that positive samples contain a superset of the information presented in negative samples that is relevant to the detection algorithm.

2) Finetuning on target and related source domain class: Mixing the source and target domains in a balanced manner during training may intuitively increase the performance of the detector on the target dataset even more, while maintaining performance on the source domain. Table IV summarizes the results of this experiment on the entire target dataset, i.e., the results are comparable to those in Table III. The 'COCO AP' column shows the AP on the person subset of COCO, to highlight performance loss on the source domain.

The 'COC+Neg./Pos./Both' entry indicates that the detector has been finetuned using the COCO 'person' subset and the negative/positive/all samples of the target dataset. First, training with negative only samples, increases the precision in comparison to both the pretrained model on COCO, as well as the finetuned models reported in Table III. Furthermore, training with positive samples significantly improves both the precision and recall scores, at the cost of 2.1% AP in the person class of COCO. Adding both positive and negative samples preserves the most pre-existing knowledge, as indicated by the small loss in person AP, as well as the 3.1% improvement in recall in the AGI dataset, in comparison to using only positive samples. These results indicate that combining source and target domain can always increase the precision compared to using data only from the target domain, as well as minimize the impact of catastrophic forgetting phenomena.

C. Qualitative Results

Figure 3 shows examples of detections made using our COCO+Both detector, on the same images as shown in Figure 1 for the pretrained model. Note that there are a lot of false positive detections using the pretrained model, which are corrected when training with the source (COCO) dataset plus the full annotated target domain dataset.

V. CONCLUSIONS

A dataset for human detection in fields was introduced, for the purposes of agricultural robotics applications. Various detection models pretrained on the VOC and COCO datasets were evaluated on this dataset, and the results indicated a severe impact of the domain shift problem. Thus, the importance of annotation of the collected images was examined, by evaluating three distinct sets of training data: a) only negative samples, i.e., no depicted humans, b) only positive samples, i.e., only images which depict humans, and c) both negative and positive. The results indicated that good performance can be achieved even when using only negative samples. However, to achieve better localization, using positive samples only is the better option. The findings of this work, along with the openly available annotated dataset, pave the way for developing methods that can work under the challenging settings that are often encountered in agricultural applications.



(a) Positive sample depicting two humans (front camera).



(b) Positive sample depicting one human (back camera).

Fig. 3: Examples of detections using our COCO+Both model: (a) Image depicting two humans, captured by the front camera of the robot, (b) Image depicting one human, captured by the back camera of the robot.

ACKNOWLEDGMENT

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 871449 (OpenDR). This publication reflects the authors views only. The European Commission is not responsible for any use that may be made of the information it contains.

REFERENCES

- M. Everingham, L. Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (VOC) challenge," *International Journal of Computer Vision*, vol. 88, pp. 303–338, 2009.
- [2] T.-Y. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2014.
- [3] N. Inoue, R. Furuta, T. Yamasaki, and K. Aizawa, "Cross-domain weakly-supervised object detection through progressive domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision* and Pattern Recognition, 2018, pp. 5001–5009.
- [4] H.-K. Hsu, C.-H. Yao, Y.-H. Tsai, W.-C. Hung, H.-Y. Tseng, M. Singh, and M.-H. Yang, "Progressive domain adaptation for object detection," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020, pp. 749–757.
- [5] X. Wei, S. Liu, Y. Xiang, Z. Duan, C. Zhao, and Y. Lu, "Incremental learning based multi-domain adaptation for object detection," *Knowledge-Based Systems*, vol. 210, p. 106420, 2020.
- [6] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *Proceedings of the European Conference on Computer Vision*. Springer, 2016, pp. 21–37.
- [7] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," arXiv preprint arXiv:1804.02767, 2018.
- [8] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang, and Q. Tian, "Centernet: Keypoint triplets for object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6569–6578.
- [9] J. P. Vasconez, G. A. Kantor, and F. A. A. Cheein, "Human-robot interaction in agriculture: A survey and current challenges," *Biosystems Engineering*, vol. 179, pp. 35–48, 2019.
- [10] T. Tommasi, M. Lanzi, P. Russo, and B. Caputo, "Learning the roots of visual domain shift," in *Proceedings of the European Conference on Computer Vision*. Springer, 2016, pp. 475–482.
- [11] S. Sankaranarayanan, Y. Balaji, A. Jain, S. N. Lim, and R. Chellappa, "Learning from synthetic data: Addressing domain shift for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3752–3761.
- [12] Y. Luo, L. Zheng, T. Guan, J. Yu, and Y. Yang, "Taking a closer look at domain shift: Category-level adversaries for semantics consistent domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2507–2516.
- [13] Y. Chen, W. Li, C. Sakaridis, D. Dai, and L. Van Gool, "Domain adaptive faster r-cnn for object detection in the wild," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3339–3348.
- [14] T. Kim, M. Jeong, S. Kim, S. Choi, and C. Kim, "Diversify and match: A domain adaptive representation learning paradigm for object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12456–12465.
- [15] R. Kemker, M. McClure, A. Abitino, T. Hayes, and C. Kanan, "Measuring catastrophic forgetting in neural networks," in *Proceedings of the* AAAI Conference on Artificial Intelligence, vol. 32, no. 1, 2018.
- [16] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama *et al.*, "Speed/accuracy trade-offs for modern convolutional object detectors," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7310–7311.
- [17] P. Nousi, E. Patsiouras, A. Tefas, and I. Pitas, "Convolutional neural networks for visual information analysis with limited computing resources," in *Proceedings of the IEEE International Conference on Image Processing (ICIP)*, 2018, pp. 321–325.