# Half a Percent of Labels is Enough: Efficient Animal Detection in UAV Imagery using Deep CNNs and Active Learning

Benjamin Kellenberger *Member, IEEE*, Diego Marcos, Sylvain Lobry, Devis Tuia, *Senior Member, IEEE*

## Abstract

**Note: this is a pre-print version of work published in IEEE Transactions on Geoscience and Remote Sensing (TGRS; in press). The paper is currently in production, the DOI link will be active soon. DOI: 10.1109/TGRS.2019.2927393.**
We present an Active Learning (AL) strategy for re-using a deep Convolutional Neural Network (CNN)-based object detector on a new dataset. This is of particular interest for wildlife conservation: given a set of images acquired with an Unmanned Aerial Vehicle (UAV) and manually labeled gound truth, our goal is to train an animal detector that can be re-used for repeated acquisitions, e.g. in follow-up years. Domain shifts between datasets typically prevent such a direct model application. We thus propose to bridge this gap using AL and introduce a new criterion called *Transfer Sampling* (TS). TS uses Optimal Transport to find corresponding regions between the source and the target datasets in the space of CNN activations. The CNN scores in the source dataset are used to rank the samples according to their likelihood of being animals, and this ranking is transferred to the target dataset. Unlike conventional AL criteria that exploit model uncertainty, TS focuses on very confident samples, thus allowing a quick retrieval of true positives in the target dataset, where positives are typically extremely rare and difficult to find by visual inspection. We extend TS with a new window cropping strategy that further accelerates sample retrieval. Our experiments show that with both strategies combined, less than half a percent of oracle-provided labels are enough to find almost 80% of the animals in challenging sets of UAV images, beating all baselines by a margin.

## Index Terms

Active Learning, Domain Adaptation, Convolutional Neural Networks, Object Detection, Unmanned Aerial Vehicles, Animal Census, Optimal Transport

## I. INTRODUCTION

R EPEATED wildlife censuses provide an invaluable tool for ecologists to count animals, monitor population health and stem threats from poaching incidents [1], [2]. Population densities and spatial locations of big mammals are constantly fluctuating, and having up-to-date information on where and how many individuals are found may be decisive for grazing needs estimation, or for the success of anti-poaching means. Hence, authorities of national parks and game reserves require animal census tools that are fast, reliable, and suitable for repeated applications over time.

Traditional censuses using manual surveys from manned helicopters [3], [4] are steadily replaced by approaches using UAVs [5]. UAVs are inexpensive, remotely-controlled aircrafts that can be equipped with small payloads like compact imaging cameras. Latest studies have shown censuses based on UAV imagery to yield superior accuracy compared to human surveys [6]. They are especially appealing when combined with methods from machine learning and computer vision [7], [8], in particular with object detectors [9], [10], [11] employing deep CNNs [12], [13]: such models allow fast scans of the thousands of images UAVs produce over game reserves of average sizes (i.e., hundreds of square kilometers), thereby alleviating the tedious work of manual photo-interpretation. This is particularly important in real-world scenarios where animals are a rare sight and images are dominated by empty background.

However, these models are typically trained on a single dataset and quickly break down in accuracy when applied to others. This problem is known as domain shift and denotes the inherent differences present between acquisitions [14]. For example, the image crops in Figure 1 are from the same game reserve, but are clearly very different in characteristics. For a human, it is trivial to locate the animals in either scene; a machine trained on only one set, however, is likely to fail when run on the other. In practice, even if it still finds most of the animals (high recall), such a model is likely to also produce false alarms everywhere in the background (low precision).

In the literature, this discrepancy is commonly solved by means of domain adaptation [14], where a model trained on one dataset (*source* domain) is modified to also work on another (*target* domain). Multiple approaches have been proposed to this end, including unsupervised ones that only consider the images of the target domain, and semi-supervised methods that further assume the presence of a small number of labels (animal positions) in the target domain.

As soon as the dominance and appearance variability of the background class get very high, a certain degree of supervision becomes unavoidable. This raises the question on how the few target labels can be obtained that are required to this end. A naive

Fig. 1: Examples from the Kuzikus dataset (see Section III-A) from 2014 (left) and 2015 (right). It is trivial for humans to identify the animals in either image, but a model trained on only one dataset is likely to fail when predicting animals in the other.

approach could require human operators to sift through hundreds of images before encountering an animal, which is highly inefficient and can lead to fatigue. This in turn likely causes erroneous labels, and hence missed targets and loss of accuracy. To this end, multiple studies have resorted to Active Learning (AL) [15]. In AL, a machine (model) works hand-in-hand with a so-called oracle (typically a human expert) and exploits their knowledge by issuing queries for ground truth whenever it encounters a particularly relevant data sample.

The notion of relevance conventionally refers to the usefulness of a sample to the final model performance on the target dataset [15]. Multiple AL criteria have been proposed [16]: for example, uncertainty sampling methods like Breaking Ties [17] exploit the model's confidence on samples; model-specific approaches like margin sampling for Support Vector Machines (SVMs) [18], expected model change [19], or the recent Bayesian CNNs [20] make use of individual model properties to establish a sample ranking. They all seek for a prioritization of samples that lead to the highest performance of the underlying model with a small, given number of queries to the oracle.

In the case of animal censuses, however, things are different: instead of improved model generalization capability, park rangers are primarily interested in *locating the animals* in the new dataset. In this context, established criteria are likely to break down, since they tend to sample in areas where the detector is uncertain and the likelihood of obtaining a true positive is very low. Finding animals thus requires an AL criterion that works in the opposite direction by prioritizing predictions that are most certainly true positives (instead of low-confidence samples).

This means that two deviating objectives need to be met: fast animal localization on the one hand, but also some model improvement on the other. The latter refers to the interactive part of the AL adaptation paradigm: one could just train a detector on source, apply it once on the target images and then use a criterion in one go. This is known as "one-shot" AL [21]. However, we argue that using the newly obtained labels at every AL iteration to update the model and re-predict candidates can lead to increasingly higher quality predictions, and thus to higher chances of true positives retrieval.

In this paper, we therefore present a novel strategy that allows finding as many animals as possible with minimal labeling effort in a new UAV acquisition, using an available source dataset and detector, AL and an oracle in the loop. In detail, the contributions of this work are as follows:

1) We introduce an AL criterion that, unlike conventional approaches, seeks to maximize the encounter probability of (rare) true positive candidates in the target domain.
2) Furthermore, we present a window cropping strategy that allows obtaining more labels per query while also being more intuitive for human annotators to label.
3) We provide an evaluation, comparison and ablation study on a UAV dataset of two distinct acquisitions characterized by domain shifts. Results show that, when using the proposed smart sampling strategy, it is possible to retrieve 80% of the animals by screening only half a percent of the acquired dataset.

The rest of this paper is organized as follows:

Section II explains the main procedure, including the AL criterion denoted as "Transfer Sampling" (Section II-A), as well as the window cropping strategy (Section II-B). We put the model to the test in Section III, results of which we show and discuss in Section IV. Finally, we draw conclusions from our work in Section V.

## II. PROPOSED METHOD

Figure 2 provides an overview of the proposed interactive domain adaptation workflow. As a precondition, it assumes the presence of a source dataset and an object detector (a deep CNN in our case) that has been trained on it. The model and its parameters are initially copied to extract features at every location in the images from the target domain. The distributions of these features in the source and target domain are then matched using Optimal Transport (OT) [22], which allows transferring the source ground truth labels to the target domain. This provides a means of confidence prediction for the target samples, which can then be verified by an expert oracle.
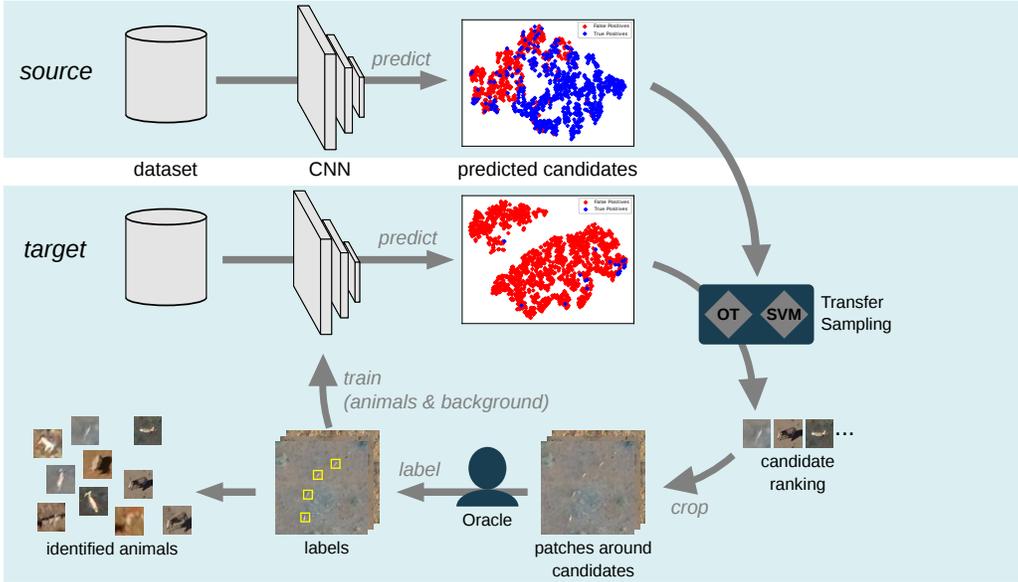


Fig. 2: Overview of the proposed workflow. We first predict candidates in the source dataset using the original, unadapted CNN (top row). We do the same on the target dataset using the current CNN (below). These serve as inputs for our TS strategy (right), which ranks the source samples with an SVM and transports the ranking to the target candidates via OT. These ranked target candidates serve as anchor points for patches, which in turn form the query data to the oracle (bottom). The latter provides labels for subsequent CNN training (completion of the AL loop).

By iteratively querying the oracle to provide ground truth labels for the most likely true positives, the model can then gradually be fine-tuned to the target domain and provides better proposals in the following AL iterations, while making sure that we minimize the tedium of the prospective oracle. We limit experiments to ground truth-based (simulated) oracles in this work, but present a further strategy to accelerate and facilitate manual annotations in upcoming extensions with human annotators, which we denote as "window cropping". In the rest of this section, we discuss the two key components of the adaptation process: the proposed AL criterion for sample selection/ranking (Section II-A), and the window cropping strategy (Section II-B).

### A. Transfer Sampling

As outlined above, our main interest lies in quickly locating animals in the target domain. Starting from the set of locations where the source model predicts more than 10%[1] chance of animal presence (denoted as *candidates* hereafter), we want to find those that are most likely to be true positives with the proposed AL criterion "Transfer Sampling" (TS).

In TS, we leverage the model's (higher) performance in the source domain and transfer this knowledge to the target samples. This is based on the assumption that the "best" predictions in source (i.e., the true animals) are clustered together in the feature space of the last layer of the CNN, and that an equivalent region can be found in the target domain that is similarly relevant.

The challenge in this context is the imbalance between animals and background, combined with a likely excessive number of false detections made by the source model in the target dataset. To find the animals quickly and keep the annotators' motivation high, we thus need to prioritize target candidates whose corresponding source predictions were indeed true positives. We therefore consider sampling according to similarities in the CNN's feature space, spanned by the deep animal detector in the source and target domains. Figure 3 shows all samples in the source domain that were predicted by the source detector as "animals". Although the model still makes a number of false positives (red), a good majority of true positives (blue) is consistently clustered in one region of the feature space. In such a scenario, it therefore makes sense to start sampling in the

---

[1]With 10% confidence we typically obtain recalls of 90% without having an excessive number of false positives.
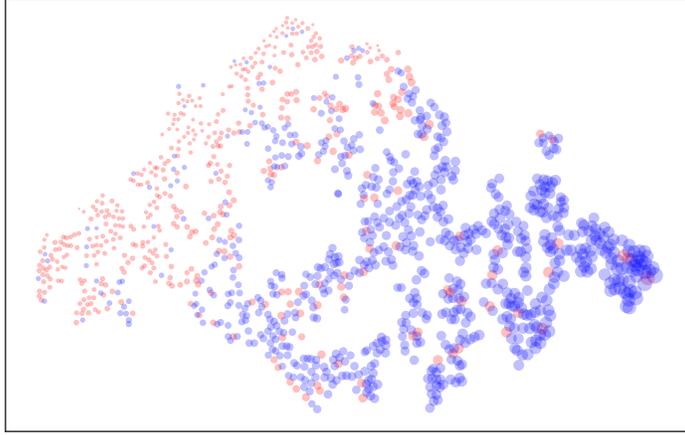
Fig. 3: Source dataset candidates with animal confidence of 0.1 or more, projected using t-SNE [23]. Blue samples were predicted correctly (true positives), red samples denote false alarms. The marker size indicates the distance to the SVM hyperplane (larger = further into the true positives region).

far right blue area, since these are the furthest away from the transition to the false positives. We thus have two tasks to solve: *(i.)* numerically identify regions in the source domain feature space that most likely contain true animals, and *(ii.)* locate the same corresponding regions in the target domain.

For the first task, we resort to a margin-based auxiliary classifier to get a surrogate measure of sample certainty. In detail, we train an SVM [24] on the full set of source candidates and then use it to rank the candidates by their distance to the separating hyperplane. This gives us an order that prioritizes samples as far away from the decision boundary as possible (the hyperplane distance is given as the marker size in Figure 3), which in turn makes sure that the most trustworthy candidates per source domain are sampled first. This strategy is conceptually close to margin sampling commonly used in AL [18], but we use it to focus the sampling on the *most* certain areas of the positive class, rather than the least certain ones.

The second task then consists of transferring the ranks to "similar" target samples. The intuition here is that both source and target candidates follow similar, mappable distributions, and we therefore need to find a way to establish an explicit source-to-target correspondence: given a predicted animal in the source domain, we want to know the predictions in the target domain that match to it. However, due to domain shifts a simple nearest neighbor search is likely to induce noise.

We propose to instead find this mapping using Optimal Transport (OT). OT finds a correspondence between two distributions that is optimal with respect to a global cost [25]. It does so by calculating and minimizing for the Wasserstein distance, also known as the earth mover's distance. This distance quantifies the difference between the two distributions as a product of their data similarities and individual distances. The intuition behind this idea is that parts of the two distributions might be similar by some measure, but far apart with respect to their "location" within the distributions. In the case of discrete distributions like ours (*i.e.*, the distributions are constituted by individual predicted animal candidates), this means that two candidates from each distribution (resp. domain) only get associated with each other if they are similar by some measure *and* lie in similar areas of their respective distributions. In the following, we therefore assume the source and target domains to be represented by the discrete probability distributions $\mu_{\mathcal{S}}$ and $\mu_{\mathcal{T}}$:

$$\mu_{\mathcal{D}} = \sum_{i=1}^{n_{\mathcal{D}}} p_i^{\mathcal{D}} \delta_{\mathbf{z}_i^{\mathcal{D}}} \text{ for } \mathcal{D} \in \{\mathcal{S}, \mathcal{T}\}, \tag{1}$$

Here, the sum over all $n$ locations (predicted candidates) of the domain $\mathcal{D}$, either source ($\mathcal{S}$) or target ($\mathcal{T}$), defines the discrete distribution. $\delta_i^{\mathcal{D}}$ denotes the Dirac at location $\mathbf{z}_i^{\mathcal{D}} \in \mathbb{R}^d$, with $\mathbf{z}_i^{\mathcal{D}}$ being the $i$th candidate's $d$-dimensional feature vector as predicted by the CNN. $p_i^{\mathcal{D}}$ is the empirical probability per sample, to which we always assign the value $p_i^{\mathcal{D}} = 1/n_{\mathcal{D}}$.

This allows us to define the OT objective for the two discrete source and target distributions: to find a set of explicit links between all the individual source and target locations that match well. To this end, OT creates a sparse matrix $\gamma$ of size $n_{\mathcal{S}} \times n_{\mathcal{T}}$, where $n_{\mathcal{S}}$ (resp. $n_{\mathcal{T}}$) is the number of samples in the source (resp. target) domain. $\gamma$ contains non-zero values wherever specific source and target locations "match". This match is defined as the link contributing to a global cost $\mathbf{C}$ for the two samples being minimal. Intuitively, establishing a link between a source and a target sample that both lie in similar
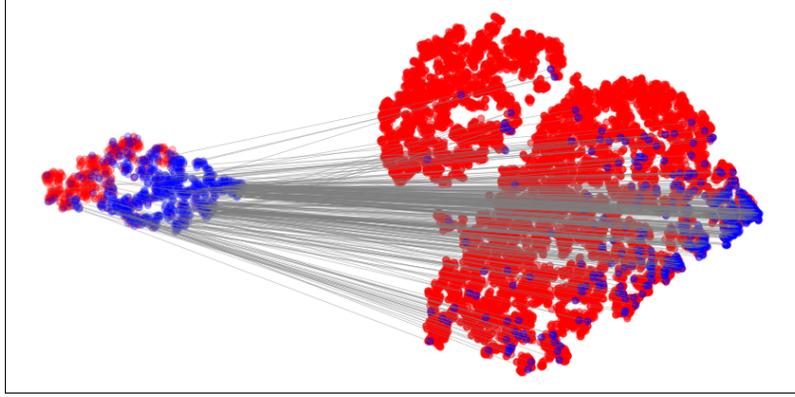
Fig. 4: A subset of predicted locations in the source (left) and target (right) domain training sets. Blue samples were predicted correctly (true positives) and red samples are false positives. Gray lines denote the correspondences found by OT for all the correctly predicted target samples. Note that, despite the imbalance and higher number of false alarms in the target set, the OT correspondences are globally consistent.

regions in the feature space induces a lower cost than if they were e.g. in opposite regions. Numerically, the optimal solution to that, i.e. the optimal *transport plan*, can be obtained as follows:

$$\gamma^* = OT(\mu_{\mathcal{S}}, \mu_{\mathcal{T}}) = \underset{\gamma \in \mathcal{B}}{\arg\min} \langle \gamma, \mathbf{C} \rangle_F, \tag{2}$$

where $\langle \cdot, \cdot \rangle_F$ is the Frobenius dot product and $\mathbf{C}$ the cost matrix of size $n_{\mathcal{S}} \times n_{\mathcal{T}}$. $\mathbf{C}_{ij}$ is the cost to move a unit amount from $z_i^{\mathcal{S}}$ to $z_j^{\mathcal{T}}$ (at source and target locations $i$ and $j$, respectively). $\mathcal{B}$ is the so-called transportation polytope, *i.e.* the set of all possible, positive matrices with prescribed marginals $\mu_{\mathcal{S}}$ and $\mu_{\mathcal{T}}$. In other words, $\mathcal{B}$ comprises all combinations of transport links between all $n_{\mathcal{S}}$ source and $n_{\mathcal{T}}$ target samples. For the cost term $\mathbf{C}$, a commonly used choice is the $\ell_2$ norm between samples [25]. We follow this approach, as we found it to work well in our setting involving CNN features. Equation (2) can be formulated as a linear program, and further be solved efficiently with simplex-based algorithms as well as group regularizations [22], [25]. This then gives us the optimal transport plan $\gamma^* \in \mathbb{R}^{n_{\mathcal{S}} \times n_{\mathcal{T}}}$, which provides explicit correspondences between individual source and target samples that are sound with respect to the whole distributions. We note that in general, and specifically when $n_{\mathcal{T}} > n_{\mathcal{S}}$, the coupling may yield one-to-many linkages, and many-to-one in the inverse case, but is always sparse thanks to the constraint that the source and target marginal probabilities (Equation (1)) must sum to one.

We can now use this transport plan to transfer the SVM-derived source scores to the individual target samples:

$$s_j^{\mathcal{T}} = \frac{1}{N} \sum_{i=1}^{n_{\mathcal{S}}} s_i^{\mathcal{S}} \delta(\gamma_{ij} > 0) \tag{3}$$

Here, $s_i^{\mathcal{S}}$ denotes the distance to the SVM hyperplane for the $i$th source sample and $\delta(\cdot)$ is the Kronecker delta, returning value 1 if the condition inside the brackets is true, and 0 otherwise. $s_j^{\mathcal{T}}$ is the score for the $j$th target sample. In essence, we assign a score to each target sample as the sum of the SVM hyperplane distances of those *source* samples whose OT link ($\gamma_{ij}$) is non-zero, normalized by $N = \sum_{i=1}^{n_{\mathcal{S}}} \delta(\gamma_{ij} > 0)$. An exemplar mapping on a subset of the training data is shown in Figure 4. This figure shows samples predicted by the CNN that has been trained on the source (left point cloud), but not yet adapted to the target domain (right point cloud). The gray lines show links obtained by OT[2]. At a first glance, it is evident that the CNN predicts orders of magnitude more false positives in the target domain, which is due to the domain shift between the two datasets. If we follow the OT links from all *source* true positives, we hit 51 target true positives (around 10% of the target true positives) and 759 target false positives (around 2.6% of the target false positives). This may sound like a low-precision result, but note that we prioritize the source true positives with our TS metric, thus drastically reducing the number of false positives (see results below). Also, the OT links to the true positive target samples consistently come from true positive source samples, which indicates that the OT-derived transport plan is globally sound and succeeds in mapping correct predictions together. In the end this means that TS is particularly robust to class imbalances: even if the ratios of true to false positives differ substantially between the two domains (as is the case in our study), TS still prioritizes the most confident predictions.

[2]For illustration purposes we only show links that point to true positives in the target domain

Once the costs are transported to the target domain, we only need to rank the target samples and can further sample them with priority on high-quality predictions.

### B. Window cropping for patch-based labeling

The second major component of our model, the window cropping strategy, extends the queried candidate with its spatial surroundings. In other words, for all candidates identified through TS, we crop a patch of fixed size around them and have the entire area labeled by the oracle, instead of the single prediction.

As mentioned above, we seek to find a trade-off between simply locating animals and CNN updates. Window cropping enhances both objectives, as it increases the total amount of labels obtainable from the oracle in a single query. We can crop a window around a query position in a UAV image in such a way that it includes as many other predicted candidates as possible. In the case of false positives, this increases the information flow to the CNN and results in it making less false predictions during the next AL iteration. However, in case neighboring candidates are also true positives, window cropping can accelerate the retrieval rate of animals with minimal additional effort from the oracle. This is not unlikely, since animals tend to flock together in groups. If the CNN thus finds just one of the animals in a herd, it is trivial for humans to localize the rest close-by this way.
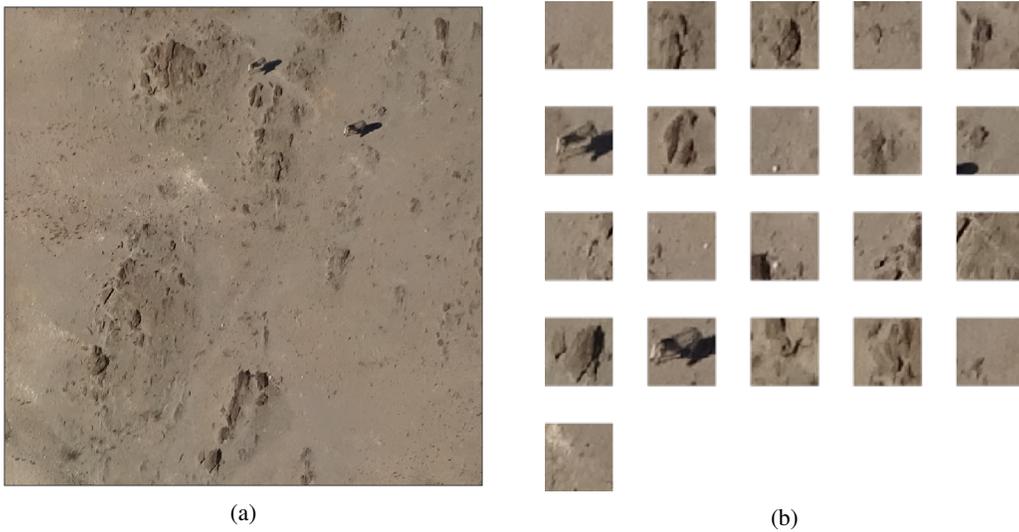


Fig. 5: Patch of a target image (left) and all candidates predicted by the source CNN in it (right). By cropping a larger patch around multiple candidates at once (left), the labeling process is both faster and more intuitive for human operators than querying on a per-candidate basis (right).

A further advantage of including the neighborhood lies in the ability of humans to be able to instantly recognize targets, if spatial context is provided. Consider Figures 5a (sample target image) and 5b (predicted candidates in it): querying the oracle for every candidate individually would not only be too exhaustive, but also more difficult, as the recognizability of the target depends heavily on spatial context, which might be missing (or confusing) on a per-sample query. In turn, locating animals in an adequately sized patch is a much simpler task for humans.

It thus makes sense to crop patches in such a way that they include as many neighboring candidates as possible. We first define a patch rectangle as $r = \{r_x, r_y, r_w, r_h\}$, with $r_x$ and $r_y$ denoting the top-left corner of it, and $r_w$ and $r_h$ the width and height in pixels, respectively. Also, let $l = \{l_x, l_y\}$ be the position in the image of the candidate selected by the AL criterion, further referred to as the *anchor point*. To select the best rectangle around the anchor point, we optimize a function that maximizes the number of candidates in the patch, minimizes the overlap with previously cropped patches and keeps the current candidate as close as possible to the center of the patch window:

$$r^* = \underset{r \in \mathcal{R}_l}{\arg\min} \left( (1 - N(p, r)) + \max(I(r, \mathcal{R}_q)) + \lambda \left\| r_c - l \right\|^2 \right) \tag{4}$$

where $\mathcal{R}_l$ is the set of windows that contain the anchor point $l$. The first term, $N(p, r)$, is the number of candidates $p$ inside rectangle $r$, normalized by the total number of candidates present in the image. The second term, $\max(I(r, \mathcal{R}_q))$, denotes the maximum area intersection between rectangle $r$ and all the rectangles in this image that have been queried before ($\mathcal{R}_q$). This term is normalized by the area of the rectangle so that it also sums to one, like the first term. The third term compares the anchor point $l$ with $r_c = \{r_x + r_w/2, r_y + r_h/2\}$, i.e. the center of the rectangle, by means of a norm and favors centering the window on the anchor. This last term primarily plays a role when the image only contains the anchor point $l$ (i.e., there

are no other candidates nor any previously queried rectangles); hence, it is downweighted with a constant $\lambda$ (set to 0.01 in the experiments). Example scenarios for the three terms are shown in Figure 6. This score function is non-differentiable in multiple ways. However, since we restrict $r$ to always contain anchor point $l$, the search space $\mathcal{R}_l$ is very limited. We thus employ an exhaustive grid search around the anchor point.
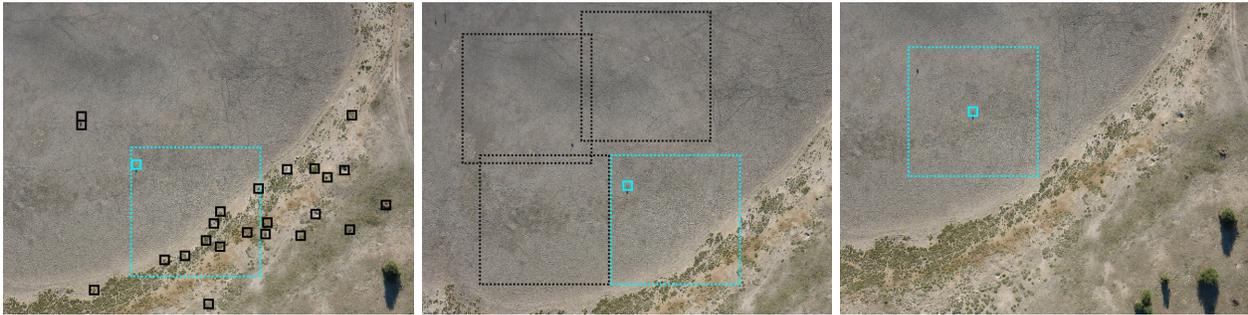


Fig. 6: For window cropping, we address different scenarios to maximize the query gain: in the first situation (left), we place the candidate window (cyan dotted) so that it includes the anchor point selected by the querying strategy (cyan solid) and as many other candidates (black) as possible. In the second situation (middle), we minimize the overlap with previously queried windows (black dotted). If neither other candidates nor previous windows are present (right), we position the window centered around the anchor point.

We then query the oracle with this patch and receive positions of animals within, if present. Any other location in the image is labeled as background.

This procedure naturally depends on the patch size, where a compromise must be found: too large patches make it increasingly harder for humans to label, while too small patches exhaust the querying budget too quickly and provide less context. In this study, we limit experiments to a simulated oracle, but nevertheless use a patch size that we found reasonable while manually labeling the dataset. In detail, we crop patches of $1000 \times 1000$ pixels (approx. $60 \times 60$m) that provide a sufficiently large number of samples while still being easy for humans to label.

## III. Experiments

We now put the proposed workflow to the test and describe the data and parameters below. Section III-A describes the two datasets used; Section III-B highlights parameters of the detector CNN and of the AL routine.

### A. Study Area and Data

We evaluate our proposed method on UAV datasets acquired over the Kuzikus wildlife reserve in Namibia[3]. Kuzikus is a private-owned park in the African savanna and home to multiple species of large mammals like kudus, giraffes, zebras, black rhinos, and more. In total, more than 3000 individuals are spread across an area of $103$km$^2$ [7], [8].

In 2014 and 2015, two image acquisition campaigns were carried out by the SAVMAP consortium[4]. A SenseFly eBee[5], equipped with a consumer-grade RGB compact digital camera, was employed for both campaigns. This resulted in 654 images for the 2014 campaign, and 3254 for the year 2015. The images of the first acquisition were initially labeled in a crowd-sourcing operation organized by MicroMappers[6] [26], followed by several iterations of refinement by the authors. The 2015 images were completely labeled by the authors.

The final statistics for both datasets are listed in Table I. Although both datasets were acquired over geographically overlapping areas, they feature a substantial domain shift in multiple ways: in terms of *external conditions*, the datasets were acquired at different times of the year (May 2014, resp. February and May 2015), under different weather and lighting conditions, with different cameras and varying flying altitudes above ground. Furthermore, additional shifts can be observed in the *label space*: the 2014 data already have a substantial class imbalance (1:10'000 in terms of animal-to-background pixels), but the 2015 dataset is larger and even more imbalanced, with an overall lower proportion of animals.

We divided the source (2014) dataset into training, validation and test splits according to the following set of rules:

- We assign entire images to only one of the three sets to avoid autocorrelation effects.
- We differentiate between images that contain at least one animal and empty ones. All the images with at least one animal are distributed so that the number of *animals* in the sets are distributed as follows: 70% for the training set, 10% for validation, and 20% for testing.

---

[3]http://kuzikus-namibia.de/xe_index.html

[4]https://lasig.epfl.ch/savmap

[5]https://www.sensefly.com

[6]https://micromappers.wordpress.com

TABLE I: Overview of the 2014 and 2015 Kuzikus UAV datasets.

|  | Set 1 | Set 2 |
|---|---|---|
| Year | 2014 | 2015 |
| Image sizes | $4000 \times 3000$ | $4896 \times 3672$, |
|  |  | $4608 \times 3456$ |
| Camera models | Canon PowerShot S110 | Sony DSC-WX220, |
|  |  | Canon IXUS 127 HS |
| No. images | 654 | 3254 |
| with animals | 239 | 111 |
| without | 415 | 3143 |
| No. animals | 1183 | 646 |
| Elevation a.g. (est.) | 120m | 160m |

TABLE II: Split properties for the 2014 and 2015 datasets.

| Set | training | | | | validation | | | | test | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | No. images | | | No. animals | No. images | | | No. animals | No. images | | | No. animals |
|  | with animals | without | total |  | with animals | without | total |  | with animals | without | total |  |
| 2014 | 159 | 291 | 450 | 830 | 35 | 41 | 76 | 118 | 45 | 83 | 128 | 235 |
| 2015 | 91 | 2750 | 2841 | 565 | 20 | 393 | 413 | 81 | - | - | - | - |

- All the remaining images (i.e., those without any animal) are then distributed at random to meet the same 70-10-20 split, but this time based on the number of *images*, as closely as possible.

For the target (2015) data, we do not require a test set, since we sample directly using the oracle. However, we do use a small validation set for hyperparameter fine-tuning. Details on the dataset splits can be found in Table II.

### B. Model Setup

In the following, we highlight the main model components and their parameters: Section III-B1 explains the deep CNN used for animal detection, and Section III-B2 provides details on the AL framework.

*1) CNN Training:* In this study, we follow the training recommendations presented in [8], which are specifically tailored to animal censuses in heavily imbalanced datasets. These recommendations include class weights; a special "border" class that is placed in the 8-neighborhood of an animal location to reduce multiple detections of it; curriculum learning, where the model is first trained on images that always contain an animal; and hard negative mining, which amplifies the weights of the four most confidently predicted false alarms after epoch 80.

We further adopt the detector CNN from [8], whose architecture is shown in Figure 7. As a feature extractor, it employs a ResNet-18 [27] that had been pre-trained on the ILSVRC dataset [28]. The model accepts image patches of size $512 \times 512$ and predicts a downsampled grid of animal probabilities ($32 \times 32$). To adapt it for predicting a grid of this size instead of a single label, the last layers, including the global average pooling layer at the end, are removed and the first convolutional layer's stride reduced to 1. We use the 512-dimensional feature vectors from the basic model (i.e., the output of the last residual block) for further adaptations. For obtaining per-location class predictions, two Multi-Layer Perceptrons (MLPs) map the feature vectors from 512 to 1024 dimensions, and then to the 3 classes (background, animal, border), respectively. Furthermore, to avoid dependencies on mini-batch configurations both during training and testing, we replace all batch normalization layers with instance normalization [29], which performs unit-norm scaling for each image in the respective mini-batch individually. We found this substitution not to harm the model performance, but to help stabilize prediction consistency.

For the fine-tuning stages throughout the AL iterations, we lower the learning rate from $10^{-6}$ (used on source) to $10^{-7}$—we found this to prevent oscillation effects on the reduced-sized target training set. Also, we disable curriculum learning and train
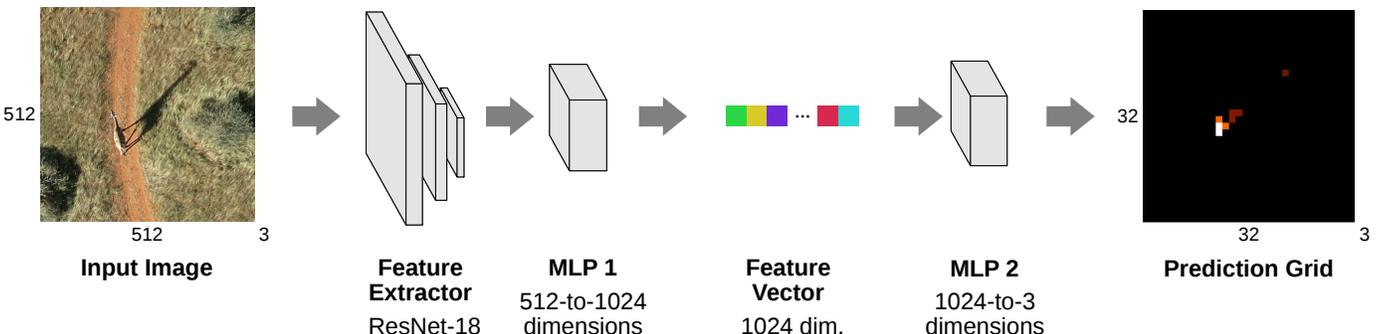


Fig. 7: Basic architecture for our animal detector, following [8]. We employ the main blocks of a ResNet-18 pretrained on ImageNet and add two more MLPs and ReLU nonlinearities. The feature vector after the MLP 1 is used by OT within TS.
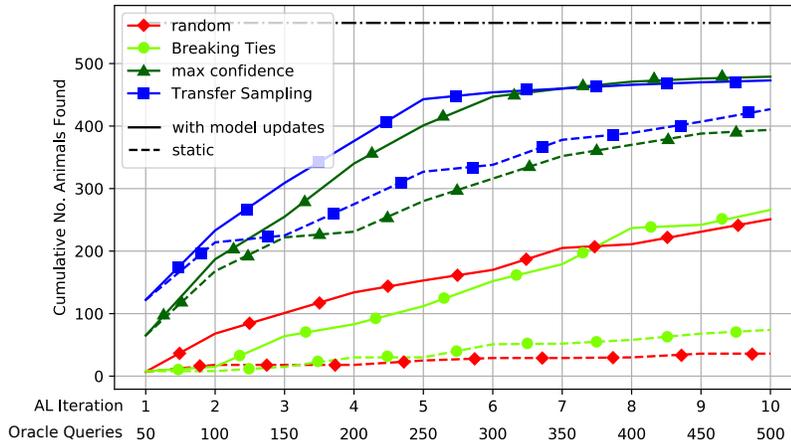
Fig. 8: Cumulative number of animals found over the AL iterations. Solid lines denote the criteria performances with model updates and target candidate re-predictions at every iteration; dashed lines are the static performances (continuous sampling on the initially predicted candidates). The black dash-dotted line marks the total number of animals in the target training set. Best viewed in color.

on the full growing AL-derived dataset at every iteration. All other parameters and training procedures, such as hard negative mining, are kept the same as in the source model.

*2) AL Loop:* In a pre-stage, we train our detector CNN until convergence on source. We use this CNN to obtain all candidates in the source training set, keeping all predictions with animal confidence of 0.1 or higher. To reduce the number of double-predictions, we employ Non-Maximum Suppression (NMS) with a search radius of 2 prediction grid cells, retaining only candidates in a 4-neighborhood with the highest confidence.

Next, we run a total of ten AL loops, querying 50 patches of $1000 \times 1000$ pixels size per iteration. We believe 50 patches to be a fair amount to query without risking the human annotators to make errors due to fatigue (note that other studies used query sizes of up to 200 images in more heterogeneous datasets [30]). Using those we fine-tune the CNN on the target training set for 12 epochs per AL iteration, which we deem a reasonable trade-off between training time and accuracy gain. Since we queried patches that are larger than what the CNN accepts as an input, we can perform extra data augmentation by randomly cropping sub-patches of the labeled areas.

We compare our TS strategy to three baselines: random candidate selection, Breaking Ties [17], and CNN confidence for being an animal ("max confidence"), all with window cropping enabled. Since model fine-tuning and candidate re-prediction is computationally intensive, we assess two scenarios for each sampling strategy: one with CNN fine-tuning at every AL iteration, and one without (i.e., using only the source model to predict candidates once and querying with TS only on those samples).

## IV. RESULTS AND DISCUSSION

Figure 8 shows the number of animals found over the course of the ten AL iterations. Already after the first 50 queries, TS found 122 animals and is significantly ahead of the baselines. This trend continues throughout the iterations, and after five AL iterations, TS found 443 out of the total 565 animals (78.4%). At this stage, the oracle had been queried 250 times. Afterwards, the total number of correctly identified animals slightly rises to 473 (83.72%). Our window cropping algorithm allows sampling patches in an out-of-grid fashion. However, if we assume uniform sampling on a grid per image, the target training set would consist of 54'324 queryable patches. This means that TS only requires the user to review around half a percent of the dataset in order to find almost 80% of the animals.

In comparison to all baselines, TS manages to yield a higher recall almost throughout the entire process. Although the max confidence ranking manages to reach roughly the same level, it does so only after the sixth AL iteration. TS in turn identified the same number of animals already an entire iteration earlier, and stayed above the rest until then by quite a margin. This means that substantially less queries need to be made to the annotator when using TS, resulting faster convergence and hence a more economical retrieval process.

Figure 9 shows the selected patches, predicted candidates and ground truth for both TS (top image) and max confidence (bottom): the latter does manage to find a reasonable number of true positives, but nonetheless misses more than half of the animals present in the scene. Explanation may be found in the t-SNE plots (right side of each image), which indicate that most of the true positives are to be found primarily in one area in the bottom right of the feature space. One might expect a correlation between feature space locations and confidences, but as shown here, this is only partially the case: TS manages
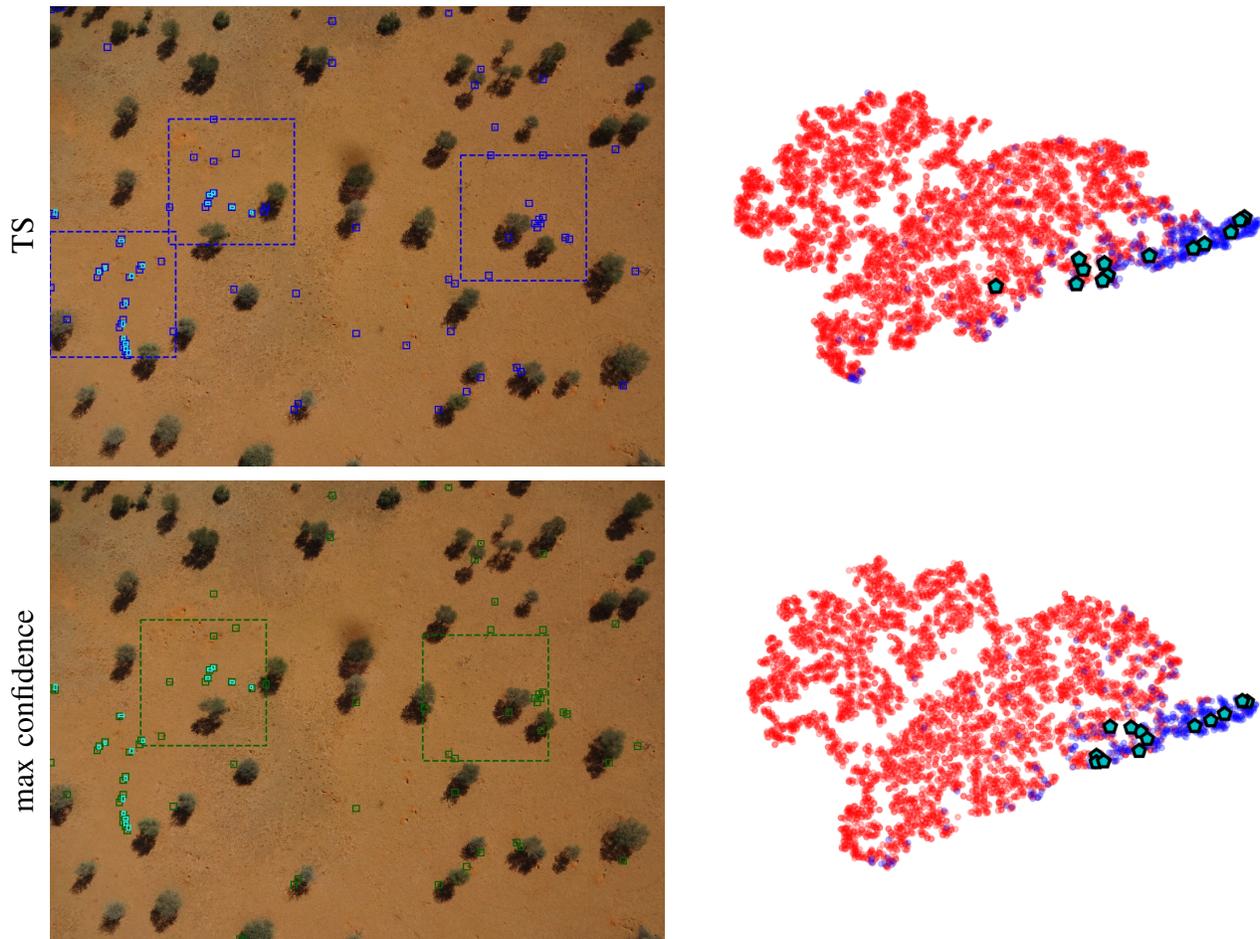
Fig. 9: Prediction examples for the TS (top) and max confidence (bottom) strategies after five AL iterations (250 queries), together with their respective t-SNE plots of the predicted candidates (right). Cyan pentagons in the t-SNE plot correspond to the cyan ground truth bounding boxes in the UAV images.

to get ahead of max confidence by strictly sampling in this area of high true positives concentration (bottom right of t-SNE plots), instead of according to confidences.

The performance of the other baselines is less satisfactory. For instance, the focus of Breaking Ties lies on minimizing the model's uncertainty. It is therefore unsurprising that it fails even worse than random sampling in finding animals. However it is still able to locate some of the animals. Consider Figure 10, which shows selected patches and predicted candidates on an example for all strategies. In this case, random sampling managed to find one of the animals in an earlier iteration, but misses the rest. Breaking Ties sampled two times in the image and found three animals, but this was most likely due to neighboring candidates that do not represent true positives (note the candidates predicted by Breaking Ties lying all on transition areas from shadow to ground). Max confidence again found the same three animals, but missed the hotspot in the center-right portion of the figure. Lastly, only TS managed to locate every single animal present, and did so with a minimal number of queries. In essence, the other strategies occasionally show true positives, but mainly due to window cropping.

Finally, in all cases we observe a significant boost when the model is updated and candidates are re-predicted versus the static "one-shot" prediction and sampling (Figure 8). This indicates that the labels provided by the oracle at each stage are useful and complete enough for adapting the CNN to the target domain, so that the predicted candidates at the next iteration are of higher quality. Window cropping helps particularly in this case, since CNNs require a large number of training samples: on the one hand, it increases the amount of image data seen by the model. On the other, it also increases the information gain: since the cropping strategy maximizes the number of predicted candidates to be included per patch, it maximizes both the number of true and false positives in the new dataset. In this respect both true and false positives are vital for the model, improving its abilities to better separate animals from background in the target dataset.

Noteworthy in this context is the difference in the number of predicted candidates between TS and the other strategies. After the fifth AL iteration, the CNN trained with TS produced a significantly higher number of candidates than all the baselines, which can be seen by the high number of blue squares in Figure 9. A possible explanation for this phenomenon is that TS finds more animals already in the first AL iteration: the CNN is then fine-tuned with a lot more true positives and therefore predicts
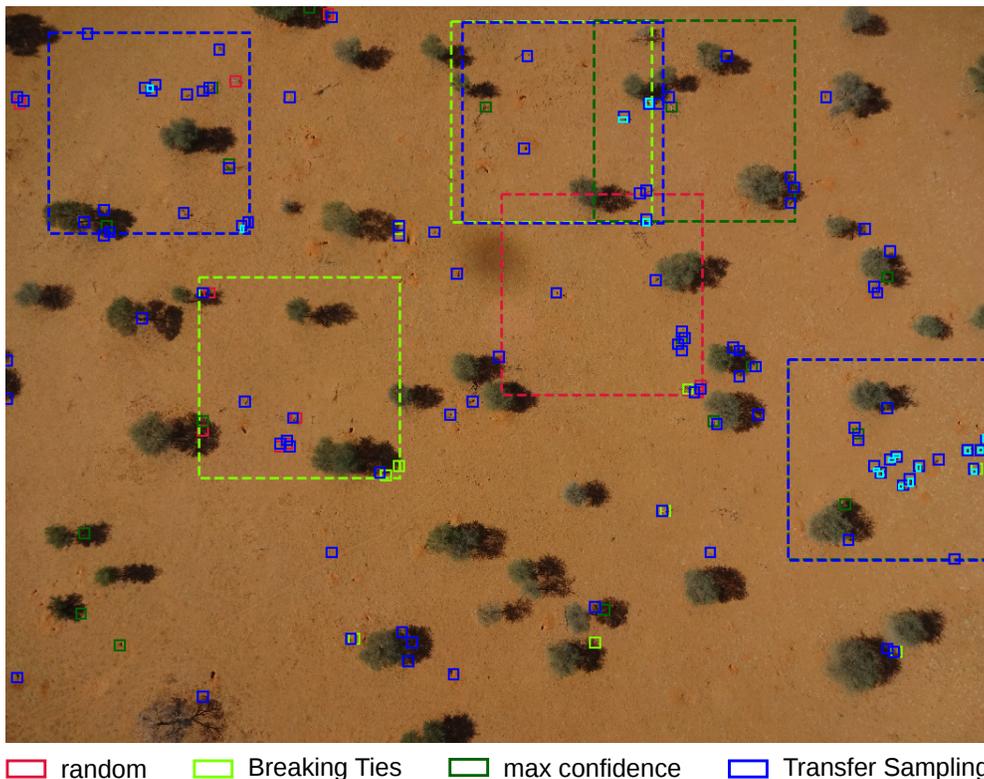
Fig. 10: Example image from the target training set with annotations after five AL iterations, showing all selected patches per sampling strategy (dashed), the candidates predicted by each CNN at the very last AL iteration, as well as the ground truth (cyan rectangles). Since only the last predictions, but all patch rectangles are shown, some of them (i.e., from earlier AL iterations) do not encompass any candidate.

more candidates when re-applied to the target dataset. In a fully-automated evaluation setting (i.e., without any oracle input), this could be problematic, since this increased number of predictions also results in more false positives. However, thanks to TS we can filter the predictions very efficiently, and as a result, localize the animals even in increasingly imbalanced settings.

## V. CONCLUSION

In this paper, we have studied the task of repeated animal censuses on UAV imagery by addressing the domain adaptation problem involved. To do so, we integrated a deep CNN-based animal detector in an AL loop. The core component of our strategy is the AL criterion: unlike traditional approaches that seek to maximize the model performance on the new dataset, our *Transfer Sampling* (TS) criterion is designed to localize the rare animals in tens of thousands of false alarms as efficiently as possible. TS works by leveraging the superior performance of the CNN detector in the source dataset (which it had been trained on) and transferring this knowledge to the target set using the distribution-mapping framework Optimal Transport (OT). The number of hits was further raised by integrating a smart window cropping strategy that maximizes the number of detections to be labeled per query, while making the labeling process itself more intuitive. Our experiments in the Namibian natural reserve Kuzikus have shown that TS indeed outperforms other AL criteria by a large margin and allows retrieving 78.4% of the animals in just 250 queries, resp. by having the oracle review less than half a percent of the entire dataset. In effect, this method thus allows for efficient and economic repetitions of animal censuses as it integrates both the adaptation and required manual verification stages into one optimized, interactive workflow.

Future work may extend this concept in multiple ways: for example, experiments with human annotators instead of simulated oracles would highlight requirements by park rangers. Extending TS e.g. by a measure of the user's confidence in providing a ground truth [31] could improve the real-world applicability of such a system. On a different track, adaptations to other geographical areas instead of new acquisitions over the same game reserve would allow testing the strategy under potentially even stronger domain shifts.

REFERENCES

[1] A. Hodgson, N. Kelly, and D. Peel, "Unmanned aerial vehicles (UAVs) for surveying Marine Fauna: A dugong case study," *PLoS One*, vol. 8, no. 11, pp. 1–15, 2013.
[2] Z. Yang, T. Wang, A. K. Skidmore, J. De Leeuw, M. Y. Said, and J. Freer, "Spotting East African mammals in open savannah from space," *PLoS One*, vol. 9, no. 12, pp. 1–16, 2014.
[3] P. Bayliss and K. Yeomans, "Distribution and abundance of feral livestock in the 'top end' of the northern territory (1985-86), and their relation to population control." *Wildlife Research*, vol. 16, no. 6, pp. 651–676, 1989.
[4] M. Norton-Griffiths, *Counting animals*. Serengeti Ecological Monitoring Programme, African Wildlife Leadership Foundation, 1978, no. 1.
[5] J. Linchant, J. Lisein, J. Semeki, P. Lejeune, and C. Vermeulen, "Are unmanned aircraft systems (UASs) the future of wildlife monitoring? A review of accomplishments and challenges," *Mammal Review*, vol. 45, no. 4, pp. 239–252, 2015.
[6] J. C. Hodgson, R. Mott, S. M. Baylis, T. T. Pham, S. Wotherspoon, A. D. Kilpatrick, R. Raja Segaran, I. Reid, A. Terauds, and L. P. Koh, "Drones count wildlife more accurately and precisely than humans," *Methods in Ecology and Evolution*, vol. 2018, no. December 2017, pp. 1–8, 2018.
[7] N. Rey, M. Volpi, S. Joost, and D. Tuia, "Detecting animals in African Savanna with UAVs and the crowds," *Remote Sensing of Environment*, vol. 200, pp. 341–351, 2017.
[8] B. Kellenberger, D. Marcos, and D. Tuia, "Detecting mammals in UAV images: Best practices to address a substantially imbalanced dataset with deep learning," *Remote Sensing of Environment*, vol. 216, pp. 139–153, 2018.
[9] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," *Advances in Neural Information Processing Systems (NIPS)*, vol. 28, pp. 1–10, 2015.
[10] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, Real-Time Object Detection," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
[11] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
[12] A. Krizhevsky, I. Sulskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," *Advances in Neural Information Processing Systems (NIPS)*, pp. 1–9, 2012.
[13] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, p. 436, 2015.
[14] D. Tuia, C. Persello, and L. Bruzzone, "Domain adaptation for the classification of remote sensing data: An overview of recent advances," *IEEE Geoscience and Remote Sensing Magazine*, vol. 4, no. 2, pp. 41–57, 2016.
[15] D. Tuia, M. Volpi, L. Copa, M. Kanevski, and J. Munoz-Mari, "A Survey of Active Learning Algorithms for Supervised Remote Sensing Image Classification," *IEEE Journal of Selected Topics in Signal Processing*, vol. 5, no. 3, pp. 606–617, 2011.
[16] B. Settles, "Active learning," *Synthesis Lectures on Artificial Intelligence and Machine Learning*, vol. 6, no. 1, pp. 1–114, 2012.
[17] T. Luo, K. Kramer, D. B. Goldgof, L. O. Hall, S. Samson, A. Remsen, and T. Hopkins, "Active learning to recognize multiple types of plankton," *Journal of Machine Learning Research (JMLR)*, vol. 6, pp. 589–613, 2005.
[18] G. Schohn and D. Cohn, "Less is more: Active learning with support vector machines," in *International Conference on Machine Learning (ICML)*, 2000, pp. 839–846.
[19] W. Cai, Y. Zhang, and J. Zhou, "Maximizing expected model change for active learning in regression," in *IEEE International Conference on Data Mining (ICDM)*, 2013, pp. 51–60.
[20] Y. Gal, R. Islam, and Z. Ghahramani, "Deep bayesian active learning with image data," *Advances in Neural Information Processing Systems (NIPS) workshop*, 2017.
[21] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap, "One-shot learning with memory-augmented neural networks," *arXiv preprint arXiv:1605.06065*, 2016.
[22] M. Cuturi, "Sinkhorn distances: Lightspeed computation of optimal transport," in *Advances in Neural Information Processing Systems (NIPS)*, 2013, pp. 2292–2300.
[23] L. J. P. Van Der Maaten and G. E. Hinton, "Visualizing high-dimensional data using t-sne," *Journal of Machine Learning Research (JMLR)*, vol. 9, pp. 2579–2605, 2008.
[24] C. Cortes and V. Vapnik, "Support vector machine," *Machine Learning (ML)*, vol. 20, no. 3, pp. 273–297, 1995.
[25] N. Courty, R. Flamary, D. Tuia, and A. Rakotomamonjy, "Optimal transport for domain adaptation," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 39, no. 9, pp. 1853–1865, 2017.
[26] F. Ofli, P. Meier, M. Imran, C. Castillo, D. Tuia, N. Rey, J. Briant, P. Millet, F. Reinhard, M. Parkan *et al.*, "Combining human computing and machine learning to make sense of big (aerial) data for disaster response," *Big Data*, vol. 4, no. 1, pp. 47–59, 2016.
[27] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 7, no. 3, pp. 171–180, 2015.
[28] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.
[29] D. Ulyanov and A. Vedaldi, "Instance Normalization: The Missing Ingredient for Fast Stylization," *arXiv preprint arXiv:1607.08022v3*, 2016.
[30] C.-C. Kao, T.-Y. Lee, P. Sen, and M.-Y. Liu, "Localization-Aware Active Learning for Object Detection," *Asian Conference on Computer Vision (ACCV)*, 2018.
[31] D. Tuia and J. Munoz-Mari, "Learning user's confidence for active learning," *IEEE Transactions on Geoscience and Remote Sensing (TGRS)*, vol. 51, no. 2, pp. 872–880, 2013.