

Deep object detection for waterbird monitoring using aerial imagery

Krish Kabra^{*,†,‡}, Alexander Xiong^{†,‡}, Wenbin Li^{†,‡}, Minxuan Luo[‡], William Lu[‡],
Tianjiao Yu[‡], Jiahui Yu[‡], Dhananjay Singh[‡], Raul Garcia[‡], Maojie Tang[‡],
Hank Arnold[§], Anna Vallery[§], Richard Gibbons[¶], Arko Barman[§]

[‡]Rice University, Houston, TX 77005, USA

[§]Houston Audubon Society, Houston, TX 77079, USA

[¶]American Bird Conservancy, The Plains, VA 20198, USA

Abstract—Monitoring of colonial waterbird nesting islands is essential to tracking waterbird population trends, which are used for evaluating ecosystem health and informing conservation management decisions. Recently, unmanned aerial vehicles, or drones, have emerged as a viable technology to precisely monitor waterbird colonies. However, manually counting waterbirds from hundreds, or potentially thousands, of aerial images is both difficult and time-consuming. In this work, we present a deep learning pipeline that can be used to precisely detect, count, and monitor waterbirds using aerial imagery collected by a commercial drone. By utilizing convolutional neural network-based object detectors, we show that we can detect 16 classes of waterbird species that are commonly found in colonial nesting islands along the Texas coast. Our experiments using Faster R-CNN and RetinaNet object detectors give mean interpolated average precision scores of 67.9% and 63.1% respectively.

Index Terms—Object detection, Convolutional neural networks, Wildlife monitoring

I. INTRODUCTION

Colonial waterbird nesting islands can be found across the globe, and each of North America’s coasts is home to its own species of breeding colonial waterbirds. Colonial waterbirds are important indicators of ecosystem health [1], provide numerous ecosystem services, and are an important part of a growing nature-based tourism sector of the economy [2]. Therefore, continuing research and monitoring of these species is critical to inform conservation decisions, encourage management of habitats for the benefit of colonial waterbirds, and to continue to gauge the surrounding ecosystem health.

Monitoring of waterbirds at colonial nesting islands is a widespread technique used to track population trends. There are many colonial waterbird monitoring programs in the U.S. including the Texas Colonial Waterbird Survey, which is one of the longest running programs that monitors waterbirds across the entire Texas coast annually since 1976 [3]. Censusing waterbirds on islands is no small task. Traditional monitoring studies of waterbirds have been conducted by traversing the colony on foot, surveying via boat, or surveying aurally using small, manned aircraft. Each of these methods has its own set of challenges and consequences. Surveying waterbirds by foot can disturb both the species of interest and the habitat

occupied. Low vantage points of boat-based surveys can result in the risk of missing nests, particularly on larger and higher islands. Moreover, in certain conditions, accessing the islands by boat can be tricky due to inclement weather conditions. Manned aerial surveys is the preferred technique by state and federal wildlife agencies, but these are expensive and require the proper conditions. In fact, aircraft crashes and boating accidents have been found to be the largest causes of mortality and injury among biologists in the field [4].

In recent years, unmanned aerial vehicles (UAVs), also referred to as drones, have presented themselves as a useful tool in wildlife management [5], [6], including waterbird monitoring [7], [8]. Drones allow researchers to remain safely on the ground while surveying areas of interest with both less cost and greater ease than traditional aerial surveys. In studies where this technology has been applied, the use of drones was found to result in more precise count estimates than traditional ground-based surveys [9]. Unfortunately, the expertise and time required to manually localize and classify species from hundreds, potentially thousands, of aerial images represents a major bottleneck.

To alleviate this issue, we developed a object detection-based deep learning pipeline that utilizes convolutional neural networks (CNNs) [10]–[12] to precisely localize and classify colonial waterbird species from UAV aerial imagery via supervised learning. We collect survey images from three colonial nesting islands along the Texas coast, and train a CNN-based object detection model to detect 16 classes of waterbirds, including the 14 most common colonial waterbird species found on these islands: Brown Pelican (*Pelecanus occidentalis*), Laughing Gull (*Leucophaeus atricilla*), Royal Tern (*Thalasseus maximus*), Sandwich Tern (*Thalasseus sandvicensis*), Great Egret (*Ardea alba*), Cattle Egret (*Bubulcus ibis*), Snowy Egret (*Egretta thula*), Reddish Egret (*Egretta rufescens*), American White Ibis (*Eudocimus albus*), Great Blue Heron (*Ardea Herodias*), Black-crowned Night Heron (*Nycticorax nycticorax*), Tri-colored Heron (*Egretta tricolor*), Roseate Spoonbill (*Platalea ajaja*), and Black Skimmer (*Rynchops niger*). We present results using two of the most commonly implemented CNN-based object detection models, Faster R-CNN [11] and RetinaNet [12].

[†]Denotes equal contribution. ^{*}Corresponding author: kk80@rice.edu

A. Contributions

The key contributions of this work are as follows:

- We develop a deep learning pipeline to detect waterbirds from UAV aerial imagery for precise waterbird monitoring. Our pipeline is general and can be applied to other applications requiring object detection from high-resolution aerial imagery, including other wildlife monitoring applications. Our code is available at: https://github.com/RiceD2KLab/Audubon_F21
- We apply our method to detect 16 classes of waterbirds from UAV aerial imagery collected from nesting islands surveyed along the Texas coast. To the best of the authors' knowledge, this is one of the largest number of species detected by a single model for UAV-based waterbird monitoring research.
- We present experimental results utilizing Faster R-CNN and RetinaNet object detectors. We show that we can accurately detect 3 of the most prevalent waterbird classes in our dataset ($> 70\%$ of total waterbirds): Mixed Tern Adults, Laughing Gull Adults, and Brown Pelican Adult, with an interpolated average precision ($AP_{IoU=0.5}$) score of over 90% for Faster R-CNN and 85% for RetinaNet. Across all waterbird classes, we achieve a mean interpolated average precision ($mAP_{IoU=0.5}$) scores of 67.9% and 63.1% for Faster R-CNN and RetinaNet respectively.

II. RELATED WORK

In recent years, object detection, the task of localizing one or more objects in an image with corresponding classifications, has seen immense advancements largely due to the rapid development of deep learning [13]. State-of-the-art object detection architectures utilizing convolutional neural networks (CNN) as a 'backbone' have particularly achieved much success due to a CNN's ability to learn hierarchical image features [14], [15]. For this work we utilize two popular CNN-based object detectors: Faster R-CNN [11] and RetinaNet [12]. Nevertheless, the proposed method is general, and is easily extensible to other object detectors.

Consequent to the success of CNN-based object detection and wide availability of open-source code, several works have utilized these methods for wildlife monitoring with unmanned aerial vehicle (UAV) imagery. Andrew *et al.* [16] use a R-CNN [10] to detect Holstein Friesian cattle from UAV imagery, proposing both a standard still-image acquisition pipeline and an extended video monitoring pipeline. Kellenberger *et al.* [17] use a custom one-stage detector with an AlexNet [14] backbone to detect large animals from UAV images captured over the Kuzikus wildlife reserve park in Namibia. Gray *et al.* [18] use a Mask R-CNN [19] to detect and segment humpback whales, minke whales, and blue whales from UAV imagery collected off the coast of California and along the Western Antarctic Peninsula.

More related to this work, researchers have also utilized CNN-based object detectors for bird monitoring. As compared to the aforementioned works, which focused on detecting

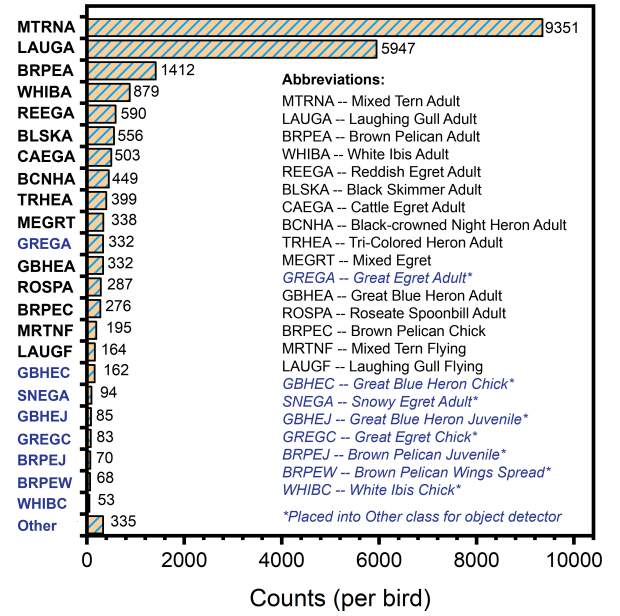


Fig. 1. **Dataset distribution of waterbird classes.** 24 classes of various bird species at different ages and configurations. In this work, we select 15 unique classes to be detected by the object detector, and combine the remaining classes into the “Other” waterbird class (highlighted by italicized, blue text).

relatively large and distinct mammals, bird detection from UAV imagery is generally regarded as a more challenging detection task due to the unique characteristics of birds. In particular, visual differences between bird species may be minor, making it difficult to distinguish between them. This difficulty is heightened for partially occluded birds, such as birds with necks tucked under their wings, as key visual features used to make distinctions are hidden. Borowicz *et al.* [20] use DetectNet [21] to count Adélie penguins in the Danger Islands off the northern tip of the Antarctic Peninsula. Hong *et al.* [22] survey various off-the-shelf CNN-based object detectors to detect birds from UAV imagery collected of both wild and decoy birds in various environments across South Korea. Hayes *et al.* [23] use a RetinaNet to detect seabirds, specifically Black-browed Albatrosses and Southern Rockhopper Penguins, from UAV imagery collected of the Falkland (Malvinas) Islands.

This work expands on the existing literature of deep learning-based object detection for bird monitoring by significantly increasing the number of bird species detected by a single object detector. The aforementioned works focus on identifying 2 or fewer bird classes that are often specific to the surveyed islands. However, this work shows that CNN-based object detectors are capable of detecting several bird species, even when visually similar or imaged in challenging viewing conditions such as dense flocks or obscuring foliage.

III. DATASET

Aerial imagery from three colonial waterbird nesting islands, Chester Island, Little Bay North Island, and North Deer

Island, was captured using a DJI Matrice 300 RTK¹ quadcopter drone with a Zenmuse P1 camera attachment². A total of 200 high-resolution images are contained in the dataset, where each image is 8192×5460 pixels in resolution. Human annotations consisting of 4 bounding-box coordinates and object classes representing different waterbirds were performed for each image. Figure 1 shows the distribution of waterbird classes present in the dataset, along with a name abbreviation list. The waterbird classes were categorized based on waterbird species, maturity, and flight. Annotations were collected in this manner, as opposed to solely variations in species type, due to the large visual differences between these classes. Note that the “Mixed Tern” and “Mixed Egret” classes do not correspond to a single waterbird species, but rather a collection of visually similar waterbird species, for example, Royal Terns (*Thalasseus maximus*) and Sandwich Terns (*Thalasseus sandvicensis*) for the “Mixed Tern” class. This was done due to human annotation difficulties in identifying the different species consistently in a large-scale manner. Finally, bird species that were either not of interest to the monitoring survey or that could not be identified by annotators were labelled as “Other”.

The dataset distribution is long-tailed, with the majority of waterbird classes dominated by Mixed Tern Adults and Laughing Gull Adults. From the original 24 classes, we focus detection efforts on 15 classes: Mixed Tern Adult, Laughing Gull Adult, Brown Pelican Adult, White Ibis Adult, Reddish Egret Adult, Black Skimmer Adult, Cattle Egret Adult, Black-crowned Night Heron Adult, Tri-colored Heron Adult, Mixed Egret, Great Blue Heron Adult, Roseate Spoonbill Adult, Brown Pelican Chick, Mixed Tern Flying and Laughing Gull Flying. The remaining 9 classes are all categorized as “Other”. Therefore, a total of 16 classes are trained for detection by object detector.

IV. METHODS

An overview of our proposed pipeline is shown in Figure 2. Here, we discuss each component of the pipeline in detail.

A. Data pre-processing

Given the large resolution of the original UAV images, we cropped the images using a sliding window tiling process, which is visually described in Figure 2. Such cropping is necessary to ensure the object detector can be trained in a computationally feasible manner. We opted not to downsample the original UAV images to minimize information loss from the original images. Cropped image tiles are of resolution 640×640 pixels, and the sliding window shifts horizontally by 400 pixels, resulting in approximately 62.5% overlap between adjacent tiles. Such overlap ensures at least one tile sees a complete waterbird that is also completely visible in the original UAV image. Consequently, we also keep edge tiles by adjusting the sliding window shift so the window remains within the bounds of the original image. Bounding box annotations for waterbirds cropped by the tiling process

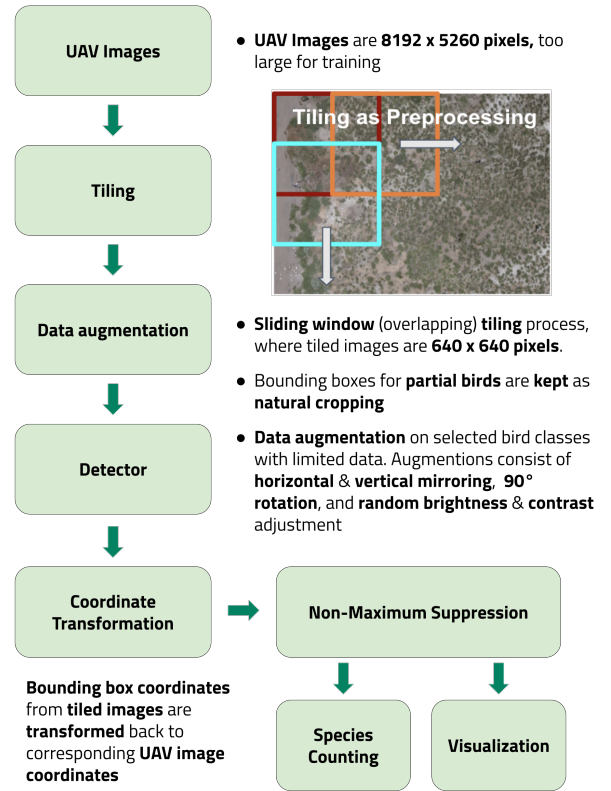


Fig. 2. **Overview of the proposed pipeline for waterbird detection from UAV aerial imagery.** The large raw UAV images are cropped using an overlapping sliding window tiling process. Visual data augmentations are performed to a subset of tiled images containing underrepresented waterbirds used for training the detector to alleviate class imbalance. During evaluation, bounding box coordinates predicted by the detector on the tiled images are transformed back into pixel coordinates in the large raw UAV image. Non-maximum suppression is used to filter overlapping detections in the final UAV image, enabling waterbird species counts to be obtained for each UAV image alongside visualizations.

are kept if there is more than 80% overlap in area with the original bounding box. Such partial birds are kept as a form of natural cropping as human annotators can also distinguish partially cropped or occluded birds.

From the tiled images, we randomly split the dataset into training, validation, and testing sets following a 70-15-15% ratio. As can be seen in Figure 1, our dataset is highly imbalanced with 3 of the classes representing over 70% of the total waterbird occurrences. To resolve this issue, we implement various image augmentation techniques by oversampling the images in which defined minority classes contain more than 80% of the total cropped image annotation. The minority classes are defined to be: Brown Pelican Adults, White Ibis Adults, Reddish Egret Adults, Tri-colored Heron Adults, Great Blue Heron Adults, Roseate Spoonbill Adults, and Brown Pelican Chicks. The augmentations consist of horizontal and vertical mirroring, 90° rotation, and random brightness and contrast adjustments.

¹<https://www.dji.com/matrice-300>

²<https://www.dji.com/zenmuse-p1>

B. CNN-based object detector

We explore two popular CNN-based object detection models, Faster R-CNN [11] and RetinaNet [12], to detect waterbirds from the tiled images. A ResNet-50 feature pyramid network [24] is utilized as the backbone for both models. For implementation, we use the Detectron2 library [25] for both models. We train both models using stochastic gradient descent with momentum ($\mu = 0.9$) [26] and a batch size of 8 for a total of 1400 steps. Both models are initialized with pre-trained COCO object detection [27] model weights. For Faster R-CNN, the initial learning rate is 0.01, which is then reduced by a factor of 0.01 after 900 steps. For RetinaNet, the initial learning rate is 0.001, which is then reduced by a factor of 0.001 after 900 steps. The learning rate and decay factor are selected using a Bayesian hyperparameter tuning algorithm [28]. Both models use the default Detectron2 L-2 weight regularization with regularization strength, $\lambda = 1 \times 10^{-4}$.

C. Data post-processing

When deploying the waterbird detector post-training, we transform the coordinates of the predicted bounding boxes in the tiled images back to the corresponding pixel coordinates of the original UAV image. Given the overlapping nature of the tiling process, there may be many detections for the same waterbird in the UAV image. To filter these overlapping bounding boxes, we utilize non-maximum suppression (NMS) with an intersection over union (IoU) threshold of 0.5 [29]. With the final predicted detections in the original UAV image coordinate space, a precise waterbird count can be made for the entire UAV image without double-counting. We note that users with access to corresponding geospatial positioning (GPS) coordinates for captured UAV images can further transform the detections to GPS coordinates. By combining detections from multiple UAV images taken during a flight mission, and filtering overlapping detections using NMS, a precise waterbird count for the area covered during the mission can be obtained.

D. Evaluation metrics

To evaluate the performance of the object detectors quantitatively, we apply the interpolated average precision (AP) metric using an IoU threshold of 0.5 and 0.75. The interpolated AP metric is commonly used by object detection competitions such as the PASCAL VOC challenge [30] and the Microsoft COCO challenge [27]. The metric defines a true positive detection to be a prediction where: (i) the confidence score of the detection is greater than a threshold, (ii) the predicted class matches the ground truth class, and (iii) the IoU between the predicted and ground truth bounding box is greater than a threshold. Here, IoU is defined mathematically as:

$$\text{IoU} = \frac{\text{Area}(B_p \cap B_{gt})}{\text{Area}(B_p \cup B_{gt})} \quad (1)$$

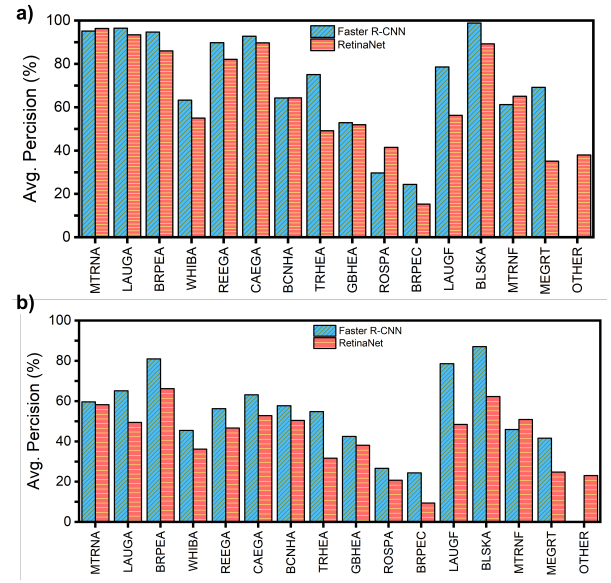


Fig. 3. **Object detector interpolated AP performance for each waterbird classes.** **A.** Bar chart for interpolated AP results using an IoU threshold of 0.5. **B.** Bar chart for interpolated AP results using an IoU threshold of 0.75.

where $B_p \cap B_{gt}$ and $B_p \cup B_{gt}$ denote the intersection and union of the predicted (B_p) and ground truth (B_{gt}) bounding boxes respectively. Precision and recall are defined as:

$$\text{Precision} = \frac{\text{True positive}}{\text{True positive} + \text{False positive}} \quad (2)$$

$$\text{Recall} = \frac{\text{True positive}}{\text{True positive} + \text{False negative}} \quad (3)$$

From Equations 1, 2, and 3, the interpolated AP is calculated as follows:

$$\text{AP}_{\text{IoU}=t} = \sum_{r \in [0:0.01:1]} (r_{i+1} - r_i) p_{\text{interp}}(r_{i+1}) \quad (4)$$

where t denotes the IoU threshold (typically ranging from 0.5 : 0.05 : 0.95), r denotes the cutoff recall, and i is the step size for recall (0.01 in our case). By calculating the AP using the formula given above, we are reducing the amount of variance “wobble” in the precision-recall curve. The interpolated precision, $p_{\text{interp}}(r)$, is defined for the maximum recall cutoff, r , and is defined as the following:

$$p_{\text{interp}}(r) = \max_{\tilde{r}: \tilde{r} \geq r} p(\tilde{r}) \quad (5)$$

In addition to interpolated AP, we construct confusion matrices to analyze the model’s classification ability between different waterbird classes that the model has already correctly localized. We define a predicted waterbird detection to be a detection with a confidence score above 0.5 and an IoU overlap between ground truth and prediction to be above 0.5. An IoU threshold of 0.5 is in line with PASCAL VOC metrics [31]. We include an additional column in the matrix to highlight missed detections by subtracting the total ground truth waterbirds to total number of predicted waterbirds that correspond with that species.

A

	Prediction																Missed
	MTRNA	LAUGA	BRPEA	WHIBA	REDEA	CAGEA	BCNHA	TRHEA	GBHEA	ROSPA	BRPEC	LAUGF	BLSKA	MTRNF	MEGRT	Other	
Ground Truth	MTRNA	1514	2	1	0	0	0	0	0	0	0	0	0	0	0	0	83
	LAUGA	3	1010	2	0	0	0	0	4	0	0	0	7	1	0	0	79
	BRPEA	0	1	251	0	0	0	0	2	0	0	0	0	0	0	0	18
	WHIBA	0	2	0	93	0	0	0	0	0	0	0	0	0	10	0	65
	REDEA	0	0	2	0	82	0	0	2	0	0	0	0	0	1	0	16
	CAGEA	0	1	0	2	0	87	0	0	0	0	0	0	0	0	0	8
	BCNHA	0	5	10	0	0	0	45	2	1	0	0	0	0	0	0	33
	TRHEA	0	13	0	0	2	0	0	39	1	0	0	0	0	0	0	16
	GBHEA	0	0	14	0	1	0	0	0	30	0	0	1	0	0	0	24
	ROSPA	0	0	1	1	1	0	0	0	6	0	0	0	0	0	0	34
	BRPEC	0	0	2	3	0	0	0	0	0	4	0	0	0	1	0	21
	LAUGF	0	5	0	0	0	0	0	0	3	0	0	15	0	1	0	10
	BLSKA	1	0	1	0	0	0	0	0	0	0	0	0	86	0	0	5
	MTRNF	4	0	0	0	0	0	0	0	0	0	2	0	11	0	0	15
	MEGRT	0	0	1	8	0	0	0	0	0	0	0	0	0	34	0	10
	Other	0	4	32	18	7	2	5	3	2	0	0	0	0	28	0	169

B

	Prediction																Missed
	MTRNA	LAUGA	BRPEA	WHIBA	REDEA	CAGEA	BCNHA	TRHEA	GBHEA	ROSPA	BRPEC	LAUGF	BLSKA	MTRNF	MEGRT	Other	
Ground Truth	MTRNA	1543	16	3	3	0	0	0	0	1	0	0	0	0	1	1	32
	LAUGA	16	1040	2	2	1	0	0	2	0	0	0	2	1	0	0	27
	BRPEA	1	2	242	0	0	0	0	1	7	1	0	2	4	0	0	11
	WHIBA	8	15	5	49	0	2	2	0	0	12	31	0	0	0	13	21
	REDEA	0	1	4	2	83	0	0	0	5	0	0	0	0	1	3	4
	CAGEA	0	1	0	14	1	73	0	0	0	1	6	0	0	0	0	2
	BCNHA	1	6	4	1	2	0	66	4	1	0	0	2	1	0	0	6
	TRHEA	0	29	1	2	10	0	2	15	8	0	0	0	0	0	1	3
	GBHEA	0	7	24	1	2	0	1	2	23	0	0	1	1	0	1	6
	ROSPA	0	1	2	0	1	0	0	0	0	30	2	0	0	0	0	5
	BRPEC	0	3	1	5	0	0	0	0	0	0	11	0	0	0	2	6
	LAUGF	0	8	5	0	0	0	1	0	1	0	0	4	0	1	1	10
	BLSKA	4	1	6	0	0	0	0	0	0	2	0	0	76	0	0	4
	MTRNF	4	1	1	0	1	0	0	0	1	0	0	0	0	11	1	11
	MEGRT	1	5	3	7	0	0	0	1	3	0	0	0	0	0	25	7
	Other	12	17	29	11	10	2	13	1	11	2	8	2	4	1	39	70

Fig. 4. **Confusion matrices for object detectors.** **A.** Confusion matrix for Faster R-CNN model. **B.** Confusion matrix for RetinaNet model.

V. EXPERIMENTAL RESULTS

The interpolated AP metrics per waterbird class for both Faster R-CNN and RetinaNet are shown in Figure 3, and the confusion matrices for Faster R-CNN and RetinaNet are shown Figure 4. We find that for both models, 6 out of 16 waterbird classes are accurately detected with high interpolated AP scores ($> 80\%$ for $\text{IoU}=0.5$, $> 60\%$ for $\text{IoU}=0.75$): Brown Pelican Adults, Laughing Gull Adults, Mixed Tern Adults, Reddish Egret Adults, Cattle Egret Adults, and Black Skimmer Adults. These 6 classes correspondingly are within the top 7 classes with the highest number of occurrences in the dataset. With regards to which waterbirds both models detect poorly, small and often hidden waterbirds, such as Brown Pelican

Chicks, Black-crowned Night Heron Adults, and Tri-colored Heron Adults, have below average interpolated AP scores. In particular, waterbird chicks are commonly found hidden behind nests, under vegetation, or under adult waterbirds, making it difficult to accurately detect such birds even by humans. Additionally, as expected, both models also perform poorly on waterbird classes with limited data such as the Roseate Spoonbill Adults and Brown Pelican Chicks, even after data augmentation, as these birds constitute less than 4% of the total population. Finally, using the confusion matrices, we find that both models become confused between waterbird classes that are visually difficult to distinguish. Both models have relatively large confusion between White Ibis Adults and Mixed Egrets, both of which consist of waterbirds with white bodies of relatively similar proportions. Similarly, Great Blue Heron Adults are commonly misclassified by the models as Brown Pelican Adults, likely due to the similar body dimensions and relative beak size when viewed from above. Lastly, we note that there is confusion between identical waterbird species that appear in different classes due to the construction of class labels to separate waterbirds in flight – namely, we see that the Laughing Gull Flying and Mixed Tern Flying classes are confused with the Laughing Gull Adult and Mixed Tern Adult classes respectively.

Overall, we find that the Faster R-CNN model outperforms RetinaNet with respect to the interpolated AP metrics, which agrees with benchmark results on standard object detection datasets and competitions. However, the confusion matrices show that, in general, Faster R-CNN has a higher false negative, or missed detection, rate than RetinaNet – specifically 14.6% over all waterbirds in comparison to RetinaNet’s 3.3%. Moreover, the confusion matrix in Figure 4A shows that Faster R-CNN fails to utilize the “Other” waterbird class at all. Finally, we find that RetinaNet can better localize waterbirds at an IoU threshold of 0.5 when compared to Faster R-CNN, but struggles to correctly classify the localized waterbird.

VI. DISCUSSION

In this work, we develop a deep learning-based tool to monitor waterbirds from UAV aerial imagery. We apply our method for precise detection of 16 classes consisting of the most common colonial waterbird species found on nesting islands along the Texas coast. To the best of the authors’ knowledge, this is one of the largest works towards combining deep learning and UAV aerial imagery for wildlife monitoring with respect to the number of species detected by a single model. We perform experiments using two popular CNN-based object detectors, Faster R-CNN [11] and RetinaNet [12]. Figure 5 visualizes sample waterbird detections on a UAV image using Faster R-CNN, illustrating the robustness at detecting waterbirds that are densely packed together or partially occluded by foliage. We quantitatively evaluate both Faster R-CNN and RetinaNet models and observe strong performance for the 3 most common waterbird classes appearing in our collected dataset: Mixed Tern Adults, Laughing Gull Adults, and Brown Pelican Adults. Specifically for Faster R-CNN, we

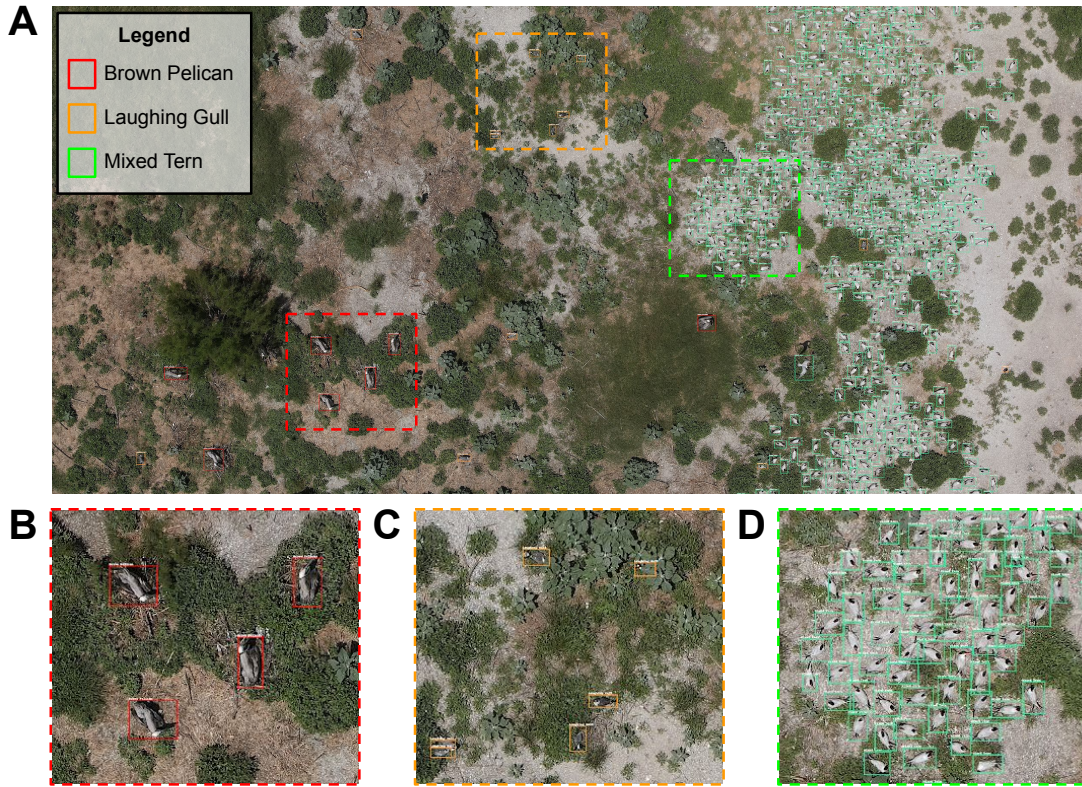


Fig. 5. Visualization of the outputs for the proposed waterbird detection from UAV aerial imagery method using Faster R-CNN. A. Original high-resolution UAV aerial image with predicted waterbird detections for 3 of the most common waterbird classes present on the surveyed islands. B. Zoomed-in image showing predictions on a flock of Brown Pelican Adults. C. Zoomed-in image showing predictions on a flock of Laughing Gull Adults. D. Zoomed-in image showing predictions on a flock of “Mixed Tern” Adults.

observe this model achieves over 90% interpolated average precision on these 3 classes.

For future work, the object detector’s relatively poor performance on minority waterbirds present in the dataset must be solved. This is a problem characteristic to learning from long-tailed data distributions and is still an active area of research [32], [33]. Furthermore, we anticipate work to further increase the number of waterbird classes detected. In particular, we emphasize the importance of separating the “Mixed” waterbird classes into their respective species, for example the “Mixed Tern” class into Royal Terns (*Thalasseus maximus*) and Sandwich Tern (*Thalasseus sandvicensis*). Such a task is extremely challenging as distinct bird species within the same family, and sometimes not within the same family, appear visually similar, even causing human experts to have trouble distinguishing species apart. This problem of fine-grained recognition is also an area of active research [34]–[36]. Finally, an analysis of the trade-off between detection precision and image resolution is necessary to understand what camera specifications and flight mission altitudes are required to monitor waterbirds adequately. This is important as care must be taken in selecting flight altitudes when surveying birds with UAVs in order to minimize disturbance, follow monitoring guidelines [37], and ensure affordable UAV and camera equipment can be utilized by surveyors.

In conclusion, we presented an artificial intelligence tool that demonstrates promising preliminary results for monitoring colonial waterbirds from UAV aerial imagery. The use of UAV aerial imagery and deep learning can advance current monitoring efforts with precise population counts in an efficient, safe, and timely manner, thereby enabling insights into more fine-scale changes in nesting habitat use year-to-year. We hope the impact of our work can help researchers and wildlife agencies alike better monitor waterbirds to gauge the ecosystem health, inform environmental policies, and protect the natural habitats of waterbirds.

ACKNOWLEDGMENTS

The project is being supported by the faculty and staff of D2K Lab at Rice University. W.L. acknowledges the National Science Foundation Graduate Research Fellowship Program. This material is based upon work supported by the National Science Foundation Graduate Research Fellowship Program under grant no. NSF 20-587. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the National Science Foundation.

REFERENCES

- [1] J. Burger, “Productivity of waterbirds in potentially impacted areas of Louisiana in 2011 following the Deepwater Horizon oil spill,”

- Environmental Monitoring and Assessment*, vol. 190, no. 3, p. 131, Feb. 2018.
- [2] U.S. Department of the Interior, U.S. Fish and Wildlife Service, and U.S. Department of Commerce, "2016 National Survey of Fishing, Hunting, & Wildlife-Associated Recreation," 2016.
 - [3] G. W. Blacklock and R. D. Slack, "The Texas Colonial Waterbird Census, 1973-1976," *Proceedings of the Colonial Waterbird Group*, vol. 2, pp. 99–104, 1979.
 - [4] D. B. Sasse, "Job-Related Mortality of Wildlife Workers in the United States, 1937-2000," *Wildlife Society Bulletin (1973-2006)*, vol. 31, no. 4, pp. 1015–1020, 2003.
 - [5] G. P. Jones IV, L. G. Pearlstine, and H. F. Percival, "An Assessment of Small Unmanned Aerial Vehicles for Wildlife Research," *Wildlife Society Bulletin*, vol. 34, no. 3, pp. 750–758, 2006.
 - [6] A. C. Watts, J. H. Perry, S. E. Smith, M. A. Burgess, B. E. Wilkinson, Z. Szantoi, P. G. Ifju, and H. F. Percival, "Small Unmanned Aircraft Systems for Low-Altitude Aerial Surveys," *The Journal of Wildlife Management*, vol. 74, no. 7, pp. 1614–1619, 2010.
 - [7] D. Chabot, S. R. Craik, and D. M. Bird, "Population Census of a Large Common Tern Colony with a Small Unmanned Aircraft," *PLOS ONE*, vol. 10, no. 4, p. e0122588, Apr. 2015.
 - [8] Y.-G. Han, S. H. Yoo, and O. Kwon, "Possibility of applying unmanned aerial vehicle (UAV) and mapping software for the monitoring of waterbirds and their habitats," *Journal of Ecology and Environment*, vol. 41, no. 1, p. 21, May 2017.
 - [9] 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. 9, no. 5, pp. 1160–1167, 2018.
 - [10] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 580–587.
 - [11] 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*, vol. 28. Curran Associates, Inc., 2015.
 - [12] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar, "Focal Loss for Dense Object Detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2980–2988.
 - [13] Z.-Q. Zhao, P. Zheng, S.-T. Xu, and X. Wu, "Object Detection With Deep Learning: A Review," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 11, pp. 3212–3232, Nov. 2019.
 - [14] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in *Advances in Neural Information Processing Systems*, vol. 25. Curran Associates, Inc., 2012.
 - [15] A. Dhillon and G. K. Verma, "Convolutional neural network: A review of models, methodologies and applications to object detection," *Progress in Artificial Intelligence*, vol. 9, no. 2, pp. 85–112, Jun. 2020.
 - [16] W. Andrew, C. Greatwood, and T. Burghardt, "Visual Localisation and Individual Identification of Holstein Friesian Cattle via Deep Learning," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2017, pp. 2850–2859.
 - [17] B. Kellenberger, M. Volpi, and D. Tuia, "Fast animal detection in UAV images using convolutional neural networks," in *2017 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, Jul. 2017, pp. 866–869.
 - [18] P. C. Gray, K. C. Bierlich, S. A. Mantell, A. S. Friedlaender, J. A. Goldbogen, and D. W. Johnston, "Drones and convolutional neural networks facilitate automated and accurate cetacean species identification and photogrammetry," *Methods in Ecology and Evolution*, vol. 10, no. 9, pp. 1490–1500, 2019.
 - [19] K. He, G. Gkioxari, P. Dollar, and R. Girshick, "Mask R-CNN," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2961–2969.
 - [20] A. Borowicz, P. McDowall, C. Youngflesh, T. Sayre-McCord, G. Clucas, R. Herman, S. Forrest, M. Rider, M. Schwaller, T. Hart, S. Jenouvrier, M. J. Polito, H. Singh, and H. J. Lynch, "Multi-modal survey of Adélie penguin mega-colonies reveals the Danger Islands as a seabird hotspot," *Scientific Reports*, vol. 8, no. 1, p. 3926, Mar. 2018.
 - [21] A. Tao, J. Barker, and S. Sarathy, "DetectNet: Deep Neural Network for Object Detection in DIGITS," Aug. 2016.
 - [22] S.-J. Hong, Y. Han, S.-Y. Kim, A.-Y. Lee, and G. Kim, "Application of Deep-Learning Methods to Bird Detection Using Unmanned Aerial Vehicle Imagery," *Sensors*, vol. 19, no. 7, p. 1651, Jan. 2019.
 - [23] M. C. Hayes, P. C. Gray, G. Harris, W. C. Sedgwick, V. D. Crawford, N. Chazal, S. Crofts, and D. W. Johnston, "Drones and deep learning produce accurate and efficient monitoring of large-scale seabird colonies," *Ornithological Applications*, vol. 123, no. 3, p. duab022, Aug. 2021.
 - [24] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.
 - [25] Y. Wu, A. Kirillov, F. Massa, W.-Y. Lo, and R. Girshick, "Detectron2," <https://github.com/facebookresearch/detectron2>, 2019.
 - [26] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *International conference on machine learning*. PMLR, 2013, pp. 1139–1147.
 - [27] 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*. Springer, 2014, pp. 740–755.
 - [28] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A next-generation hyperparameter optimization framework," *CoRR*, vol. abs/1907.10902, 2019. [Online]. Available: <http://arxiv.org/abs/1907.10902>
 - [29] T. Malisiewicz, A. Gupta, and A. A. Efros, "Ensemble of exemplar-svms for object detection and beyond," in *ICCV*, 2011.
 - [30] 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, Jun. 2010.
 - [31] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes challenge: A retrospective," vol. 111, no. 1, pp. 98–136, 2015. [Online]. Available: <https://doi.org/10.1007/s11263-014-0733-5>
 - [32] Y.-X. Wang, D. Ramanan, and M. Hebert, "Learning to Model the Tail," in *Advances in Neural Information Processing Systems*, vol. 30. Curran Associates, Inc., 2017.
 - [33] B. Kang, S. Xie, M. Rohrbach, Z. Yan, A. Gordo, J. Feng, and Y. Kalantidis, "Decoupling representation and classifier for long-tailed recognition," in *International Conference on Learning Representations*, 2020. [Online]. Available: <https://openreview.net/forum?id=r1gRTCvFvB>
 - [34] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, "The caltech-ucsd birds-200-2011 dataset," California Institute of Technology, Tech. Rep., 2011.
 - [35] S. Branson, G. Van Horn, S. Belongie, and P. Perona, "Bird species categorization using pose normalized deep convolutional nets," *arXiv preprint arXiv:1406.2952*, 2014.
 - [36] J. Krause, H. Jin, J. Yang, and L. Fei-Fei, "Fine-grained recognition without part annotations," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
 - [37] E. Vas, A. Lescroël, O. Duriez, G. Boguszewski, and D. Grémillet, "Approaching birds with drones: First experiments and ethical guidelines," *Biology Letters*, vol. 11, no. 2, p. 20140754, Feb. 2015.