

# Modeling Semantic Correlation and Hierarchy for Real-world Wildlife Recognition

Dong-Jin Kim<sup>1</sup>, Zhongqi Miao<sup>2</sup>, Yunhui Guo<sup>3</sup>, Stella X. Yu<sup>2,4</sup>, Kyle Landolt<sup>5</sup>, Mark Koneff<sup>6</sup>, Travis Harrison<sup>5</sup>

<sup>1</sup>Hanyang University <sup>2</sup>UC Berkeley / ICSI <sup>3</sup>University of Texas at Dallas <sup>4</sup>University of Michigan <sup>5</sup>United States Geological Survey <sup>6</sup>US Fish and Wildlife Service



## Abstract

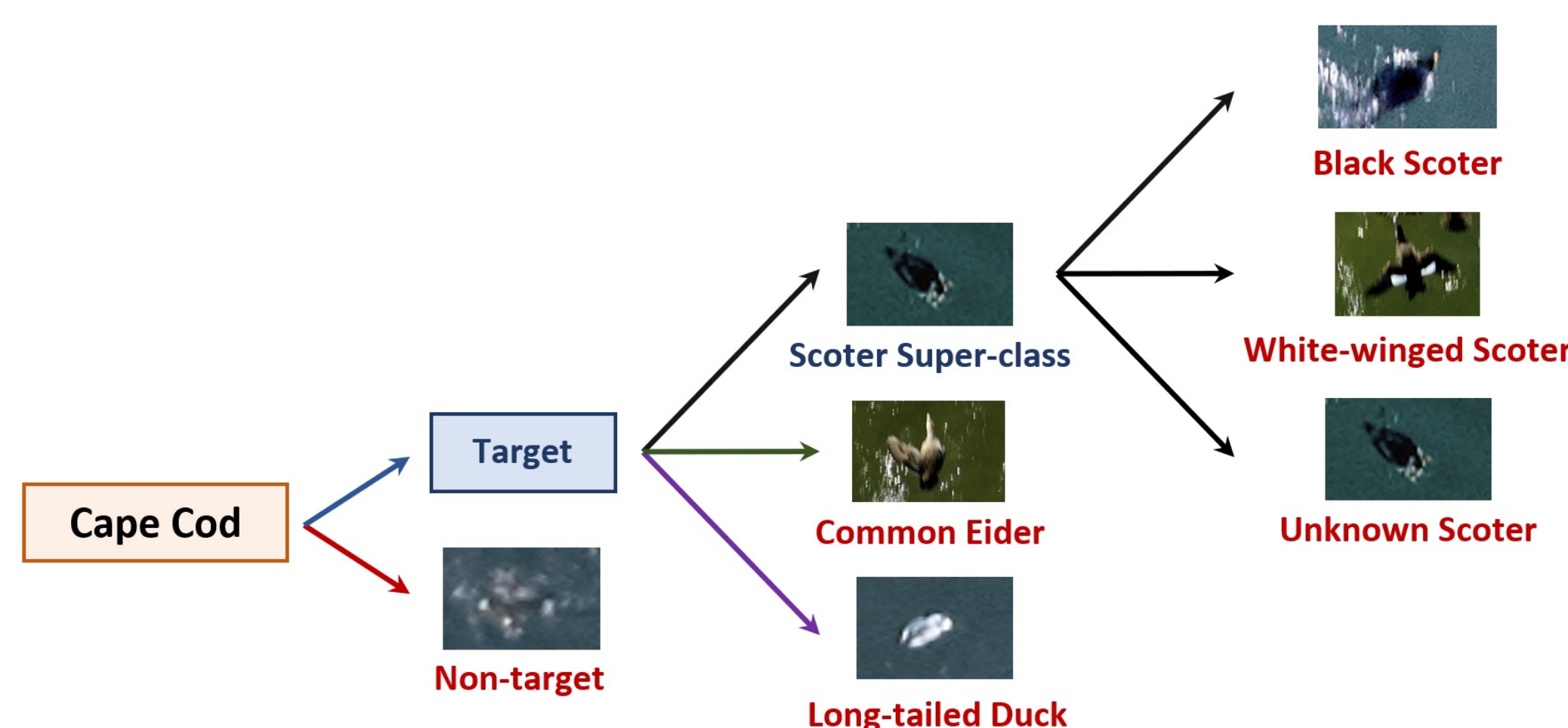
We explore the challenges of human-in-the-loop frameworks to label wildlife recognition datasets with a neural network. In wildlife imagery, the main challenges for a model to assist human annotation are two-fold: (1) the training dataset is usually imbalanced, which makes the model's suggestion biased, and (2) there are complex taxonomies in the classes. We establish a simple and efficient baseline, including the debiasing loss function and the hyperbolic network architecture, to address these issues. Moreover, we propose leveraging the semantic correlation to train the model more effectively by adding a co-occurrence layer to our model during training. We demonstrate the efficacy of our method in both our real-world wildlife areal survey recognition dataset and the public image classification dataset, CIFAR100-LT and CIFAR10-LT.

## Introduction

- Processing **areal remote sensing** datasets is expensive.
- **Human-in-the-loop** frameworks to label **wildlife** datasets.

### Main challenges:

1. Class **imbalance**. → Debiasing Loss Function
2. Complex **taxonomies**. → Hyperbolic Module & Co-occurrence



The 6 classes in our dataset have a **hierarchical** relationship [1].

## Modeling Semantic Correlation and Hierarchy

Our model consists of 3 components.

### 1. Debiasing Loss Function: Logit Adjustment [1]

$$\mathcal{L}_{LA} = -\log \frac{\exp(f_y(x) + \tau \log \pi_y)}{\sum_{j \in [L]} \exp(f_j(x) + \tau \log \pi_j)}$$

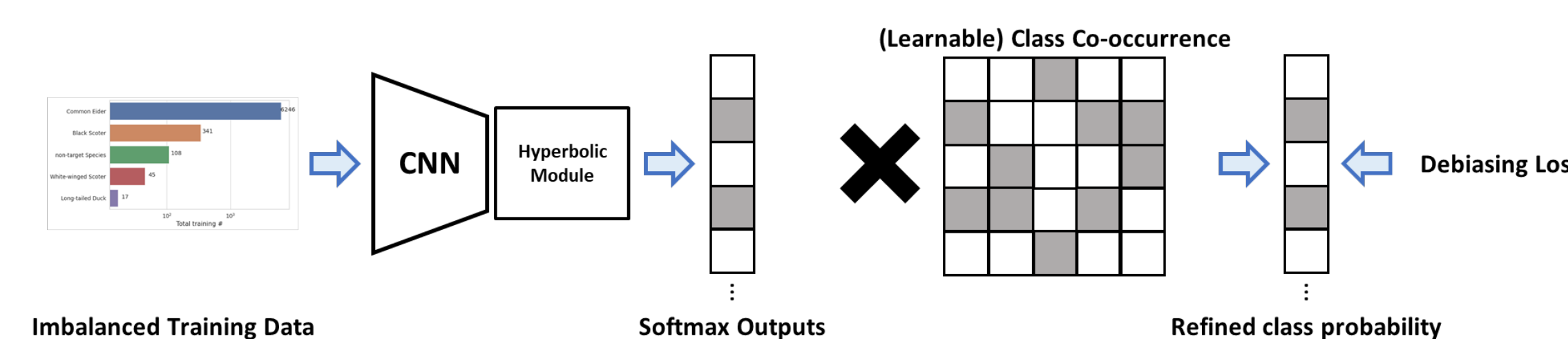
### 2. Hyperbolic Network Architecture [3]

$$\text{CLIP}(x^E; r) = \min\left\{1, \frac{r}{\|x^E\|}\right\} \cdot x^E$$

### 3. New Method: Semantic Correlations [4,5]

- A **co-occurrence matrix** to refine the class probability:

$$P(Y = i|X) = \sum_{j \in \mathcal{Y}} p(Y = i|Y' = j) * p(Y' = j|X).$$



Our framework including a **hyperbolic module**, debiasing **loss function**, and our **co-occurrence layer**.

## Experiments

Species	Train #	Test #
Common Eider	6,246	3,172
Unknown Scoter	466	114
Black Scoter	341	108
White-winged Scoter	45	21
Long-tailed Duck	17	5
Non-target Species	108	38

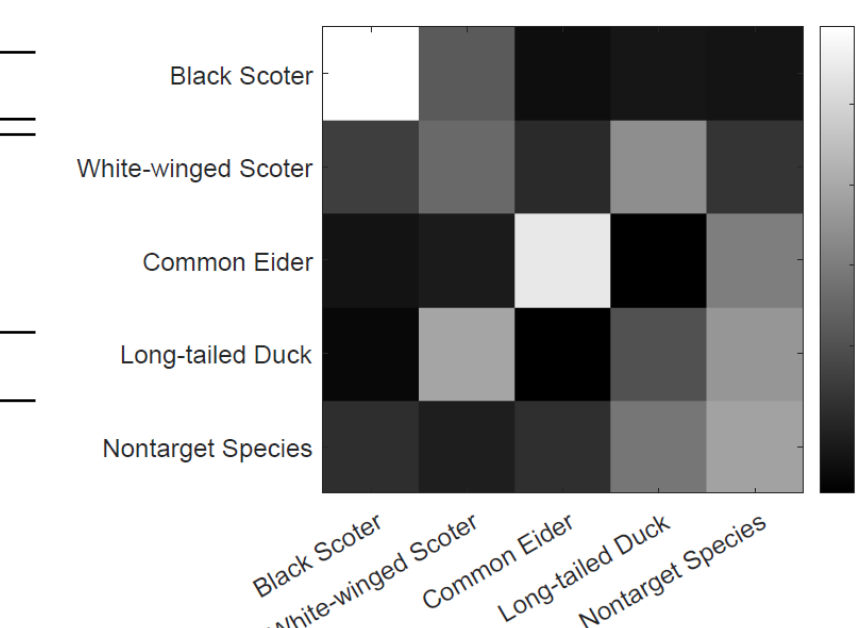
The number of samples in our dataset.

Species	Test accuracy (%)				
	Vanilla	+LDAM	+LA	+LA+H	+LA+H+Corr (Ours)
Common Eider	<b>99.1</b>	93.8	92.0	86.8	91.5
Black Scoter	95.4	92.6	97.2	87.0	<b>97.2</b>
White-winged Scoter	38.1	71.4	85.7	61.9	<b>85.7</b>
Long-tailed Duck	0.0	100.0	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>
Non-target Species	39.5	68.4	78.9	<b>84.1</b>	81.6
<b>Average accuracy (%)</b>	54.1	85.3	90.8	82.9	<b>91.5</b>

Our model with logit adjustment loss (+LA), hyperbolic module (+H), and our semantic co-occurrence method (+Corr) shows the best performance.

Methods	CIFAR100-LT	CIFAR10-LT
Vanilla	38.3	70.4
LDAM [1]	42.0	77.0
LA [16]	43.9	77.7
<b>Ours</b>	<b>45.3</b>	<b>79.0</b>

Accuracy on CIFAR10-LT and CIFAR100-LT dataset.



The visualization of the co-occurrence matrix.

## Acknowledgement

The findings and conclusions in this article are those of the author(s) and do not necessarily represent the views of the U.S. Fish and Wildlife Service.

## References

- [1] Zhongqi Miao et al., Challenges and solutions for automated avian recognition in aerial imagery. *Remote Sensing in Ecology and Conservation* 2022.
- [2] Aditya Krishna Menon et al., Long-tail learning via logit adjustment. *ICLR* 2021.
- [3] Yunhui Guo et al., Clipped hyperbolic classifiers are super-hyperbolic classifiers. *CVPR* 2022.
- [4] Dong-Jin Kim et al., Detecting human-object interactions with action co-occurrence priors. *ECCV* 2022.
- [5] Dong-Jin Kim et al., Acp++: Action co-occurrence priors for human-object interaction detection. *IEEE TIP* 2021