

Building a Flexible Framework for Automated White Shark Re-Identification

Fabrice Kurmann¹, Connor Pryor¹, Charles Dickens¹, Alexandra E. DiGiacomo², Samantha Andrzejaczek², Eriq Augustine¹, Barbara A. Block², Lise Getoor¹ University of California Santa Cruz¹, Stanford University²



Introduction and Motivation

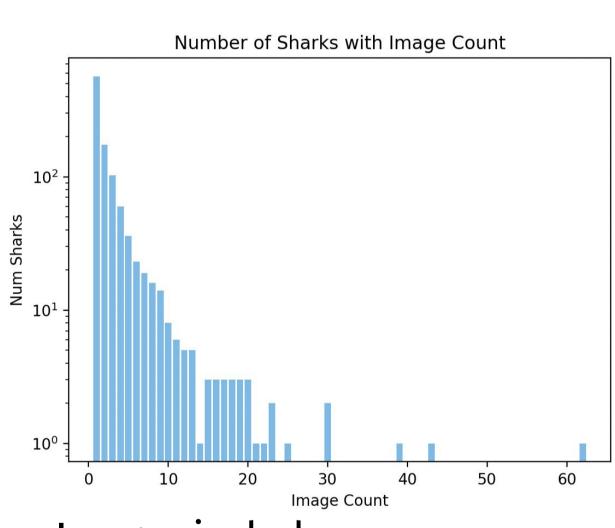
Animal Re-Identification (Re-ID), matching new observations against a catalog of known individuals, is essential for wildlife conservation, however manual re-identification is time-consuming and error-prone.

We introduce an automated white shark Re-ID framework designed to accelerate and improve this process.

Our system leverages visual features of dorsal fins [1] while keeping humans in the loop for validation, enabling efficiency and reliability.

Dataset

- Their rarity, vast habitat, and difficulty to photograph creates a sparsely populated database of many sharks, each represented by few images
- 3083 dorsal fin images
- 1031 unique white sharks
- 20+ year timespan



- Images include numerous confounding factors
- Fin orientation, rotation, angle
- Background color/lighting
- Water surface, reflection, splash



Methods

Image encoder maps images to embedding space

- Google/vit-large-patch16-384 feature extraction backbone [2]
- Multilayer perceptron (MLP) projection head
- Trained for:
- Sensitivity to biomarkers on the dorsal fin (notches and pigmentation)
- Invariance on confounding variables (pose, lighting, image quality)

Retrieval algorithm searches embeddings space

- Well-trained model organizes embeddings into distinct clusters for each shark
- Matching sharks are retrieved by searching the embedding space for nearest neighbors to a query image

Shark Re-Identification Framework Shark Catalog Image Encoder Database Rankings 1. ID 472

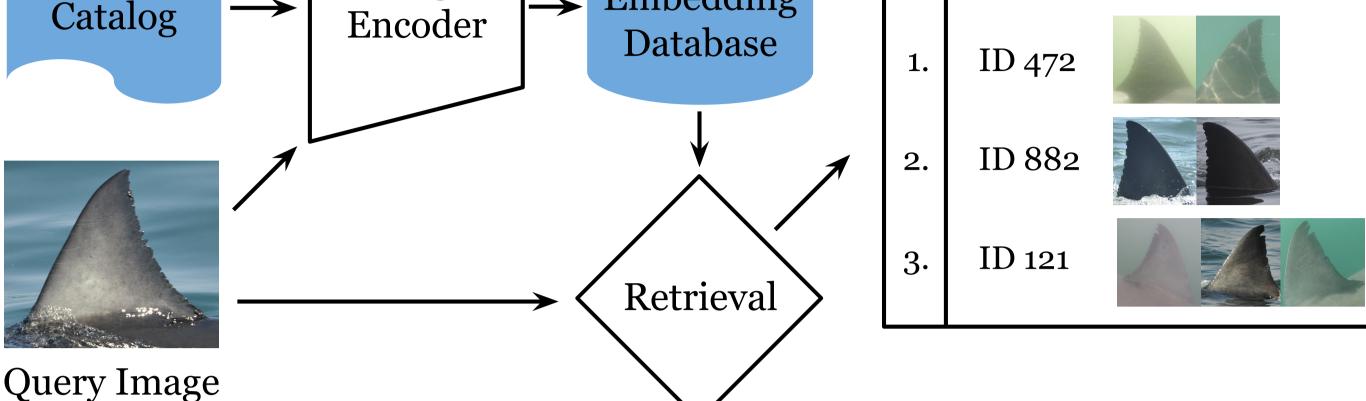


Image Encoder Training

Training Approach:

Triplet loss function [3]:

$$\mathcal{L}(A, P, N) = \max(d(A, P) - d(A, N) + \alpha, 0)$$

- Fine-tuning with low rank adaptation [4]
 - Parameter efficient method allows convergence in 72 hours on single
 GPU

Training Enhancements:

- Class-aware triplet sampling
 - Sparse dataset: many training batches lack anchor-positive samples
 - Explicitly sample two positive samples for each anchor image
- Training image augmentation
 - Augmentations to perspective, rotation, scale, and color during training
 - Increase invariance to confounding traits in images
 - Reduce overfitting

Retrieval Techniques

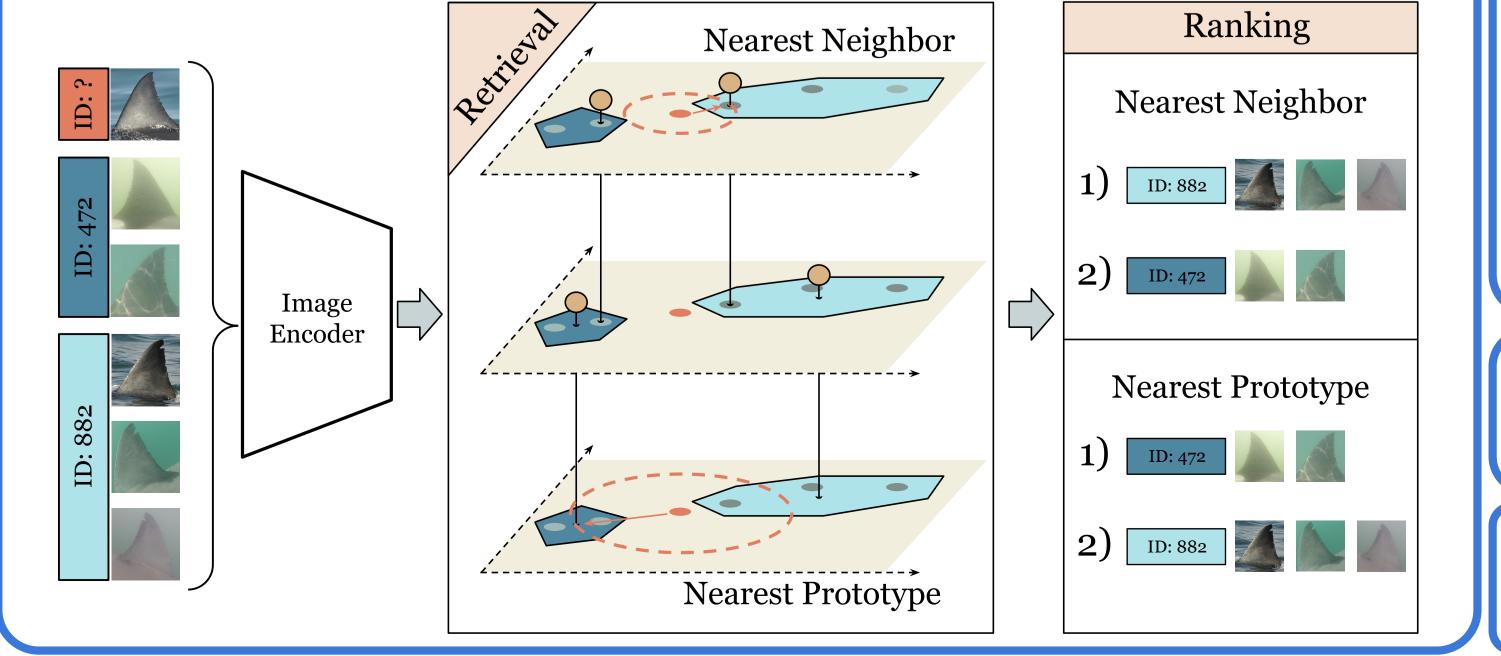
Nearest Neighbor (NN)

 Rank all identities by their closest image embedding to the query

$$\hat{\mathcal{Y}}_{ ext{NN}} = ext{rank} \left(y \in \mathcal{Y} \; \middle| \; \min_{z \in Z_y} d(z_q, z)
ight)$$

Nearest Prototype (NP)

- Prototype embedding $\mu_y = \frac{1}{|Z_y|} \sum_{z \in Z_y} z_z$ for each ID [5]
- Rank by the closest prototype to the $\hat{\mathcal{Y}}_{NP} = \operatorname{rank} \left(y \in \mathcal{Y} \,\middle|\, d(z_q, \mu_y) \right)$ query



Results

- Hits@K scores
- Proportion of queries where correct individual is among first k retrieved
- K=50 represents a practical upper limit for human review
- Test set of 1 randomly selected image from each individual with
 1 image
 - 497 train, 2,586 test images
- Ablation shows breakdown of augmentation technique.
- All techniques aside from random erasing improve results

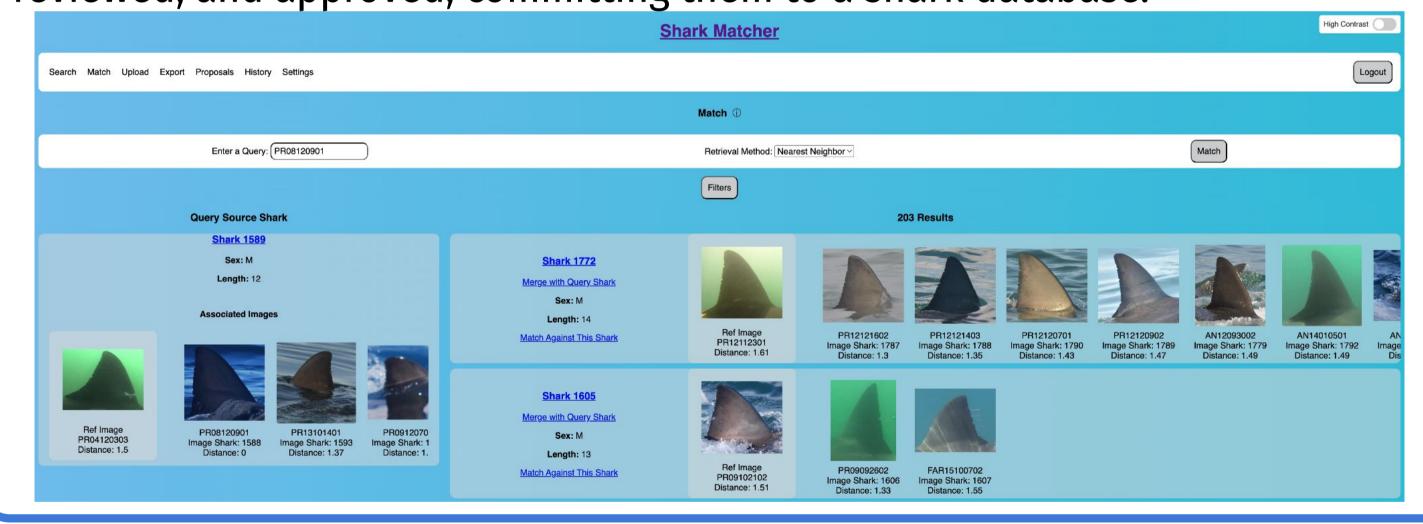
| Model | Nearest Neighbor Retrieval | | | |
|------------------------|----------------------------|--------|---------|---------|
| | Hits@1 | Hits@5 | Hits@25 | Hits@50 |
| CNN | 0.02 | 0.05 | 0.14 | 0.20 |
| ViT | 0.03 | 0.09 | 0.23 | 0.33 |
| ViT + LoRA | 0.09 | 0.23 | 0.40 | 0.53 |
| ViT + LoRA + CAS | 0.31 | 0.55 | 0.74 | 0.83 |
| ViT + LoRA + CAS + Aug | 0.48 | 0.67 | 0.85 | 0.90 |

| odel | Prototype Retrieval | | | | |
|-----------------------------|----------------------------|--------|---------|---------|--|
| | Hits@1 | Hits@5 | Hits@25 | Hits@50 | |
| IN | 0.02 | 0.07 | 0.15 | 0.22 | |
| Γ | 0.05 | 0.13 | 0.30 | 0.38 | |
| Γ + LoRA | 0.11 | 0.23 | 0.43 | 0.53 | |
| Γ + LoRA + CAS | 0.33 | 0.56 | 0.75 | 0.82 | |
| Γ + LoRA + CAS + Aug | 0.51 | 0.68 | 0.84 | 0.91 | |

| Augmentations | Nearest Neighbor Retrieval Hits@1 Hits@5 Hits@25 Hits@50 | | | |
|----------------------|--|------|------|------|
| None | 0.31 | 0.55 | 0.74 | 0.83 |
| Geometric | 0.45 | 0.65 | 0.83 | 0.89 |
| Geo. + Color | 0.48 | 0.67 | 0.85 | 0.90 |
| Geo. + Color + Erase | 0.35 | 0.57 | 0.74 | 0.83 |

Shark Matcher Human-in-the-Loop UI

Collaborating with marine biologists, we developed a UI taylored for their labeling workflow. Our model's match results can be browsed, filtered, reviewed, and approved, committing them to a shark database.



Conclusion and Future Work

We develop a framework for white shark Re-ID, presenting optimizations to model training and retrieval technique, showing their benefits to retrieval accuracy. Paired with our UI, this work has become a valuable tool for dataset de-duplication and matching newly captured shark images for our marine biologist collaborators.

As future work, we aim to:

- Evaluate over different animal species to better understand generalization
- Improve model adaptation with online learning of newly added data
- Develop techniques for training accurate models on noisy and mislabeled real-world data

Acknowledgements

This work was partially supported by the NSF grant CCF-2023495.

[1] Nowacek, Christiansen, Bejder, Goldbogen, Friedlaender. Studying cetacean behaviour: new technological approaches and conservation apps. (2016).

[2] Dosovitskiy, Beyer, Kolesnikov, Weissenborn, Zhai, Unterthiner, Dehghani, Minderer, Heigold, Gelly, Uszkoreit, and Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv (2021)

[3] Schroff, Kalenichenko, and Philbin. Facenet: A unified embedding for face recognition and clustering. CVPR (2015).

[4] Hu, Shen, Wallis, Zhu, Li, Wang, Wang, and Chen, Lora: Low-rank adaptation of large language models, arXiv (2021).