# PIVOT: CNN-Labeled Image Validation

Chandler Ault, Yamina Katariya, Yash Manne, Aditi Shrivastava, Alison Chase*, Claire Berschauer*

## Abstract

There is an urgent need to accurately quantify oceanic phytoplankton to understand their role in the global Carbon Cycle. The prevailing method is a **complex, time-consuming, scattered** process of model development, image labeling, and model evaluation. In collaboration with researchers at the University of Washington Applied Physics Lab*, we developed a **unified, interactive tool to validate the Convolutional Neural Network (CNN) classification** of phytoplankton imagery.

Our in-house solution, the **Plankton Image Validation Optimization Toolkit (PIVOT)**, enables researchers to efficiently correct mislabeled images, improve the model, and track its performance. With **active learning**, PIVOT **strategically selects training images** to enhance model performance with minimal labeling while also incorporating a subset of random images for validation. Thus, we facilitate a **human-in-the-loop approach** to evaluate model performance while also accumulating a valuable dataset for subsequent use in semi-supervised learning.

## Methods

### Exploratory Analysis

The images utilized in this project were sourced from an Imaging FlowCytobot (IFCB) during four research cruises spanning 103 different dates and 5,299 sampling points. This effort resulted in a dataset comprising approximately 5,514,006 images of single-cell phytoplankton. Subsequently, a CNN was trained to classify these images into 10 groups: **Chlorophytes**, **Ciliates**, **Cryptophytes**, **Diatoms**, **Dictyo**, **Dinoflagellates**, **Euglenoids**, **Prymnesiophytes**, **Unidentifiable**, and **Other**.
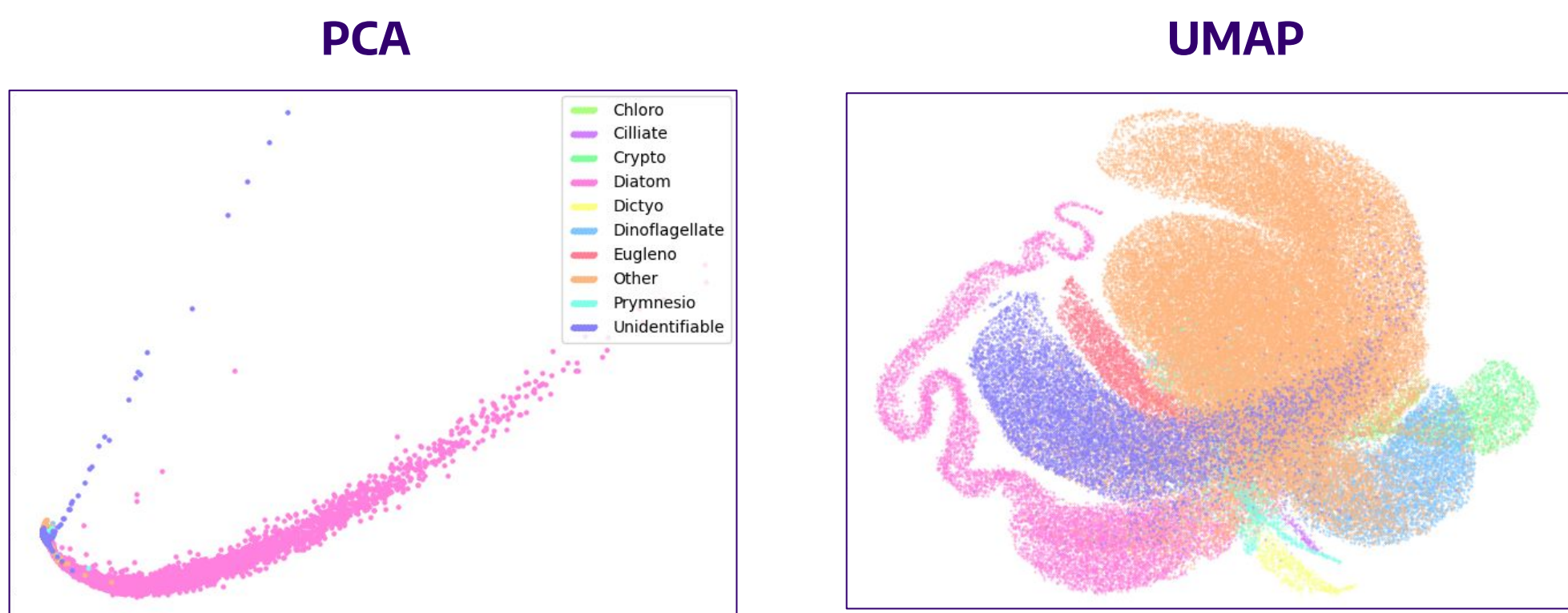
### PCA          ### UMAP



**Figure 1**: Skewed Distribution of Model Classes

In **Figure 1**, we see that most classes are fairly separable while the Unidentifiable and Other classes dominate the vector-space, potentially leading to erroneous misclassification. So, it might prove valuable to pre-filter these classes before proceeding with standard classification.

### Application

We employ **Streamlit** to craft an interactive web tool meticulously tailored for the dynamic validation of CNN-based image classifiers. Streamlit **seamlessly integrates with Python**, enabling swift development of this application with a **focus on user-friendliness and accessibility** for researchers of varying technical proficiencies. Our web tool presents an intuitive interface, empowering researchers to **label images, review summary statistics, and refine the model effortlessly**, all while preserving user progress and offering supplemental resources as needed.
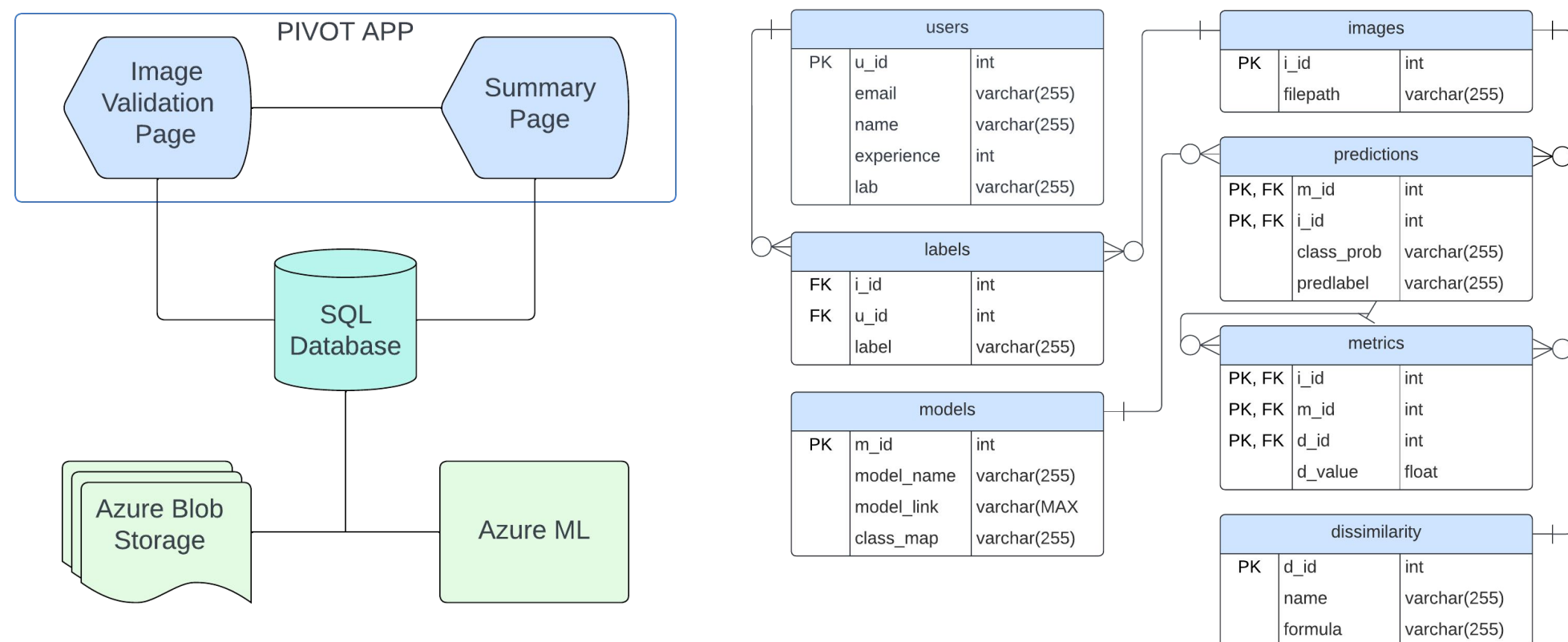


**Figure 2**: System and Entity Relationship Diagrams

### Data Storage

In our system architecture, data components are organized and stored within a **Azure SQL database** as seen in **Figure 2**. We developed utility functions and stored procedures that **dynamically retrieve data** based on user interactions and system events, such as model updates or new image availability. This database serves as the backbone of our application, enabling seamless data retrieval and manipulation for display on the website.

### Model Pipeline

The deployment and hosting of the pre-trained CNN model are managed through **Azure Machine Learning**. We find that this service **prioritizes efficient deployment** and allows researchers to use the model without extensive setup. We conduct the entire process shown in **Figure 3** on the Azure virtual machine provided by the UTOPIA group at the eScience Institute.
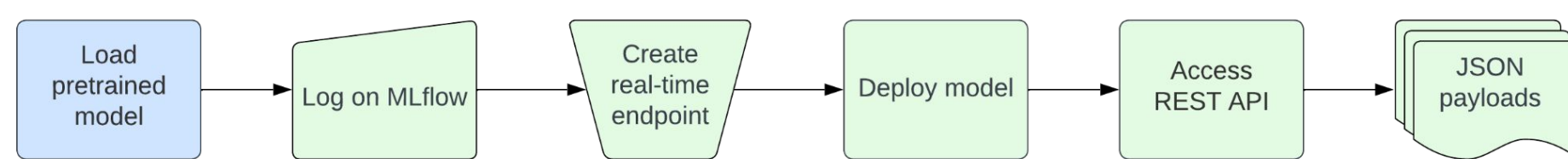


**Figure 3**: Azure Machine Learning Flow Diagram

### Image Recommendation Optimization

We implement a framework to precompute an image importance score for each image per model. The **top N images** are then selected for validation for each user session. Additionally, a random selection of **R hold-out evaluation images** are also selected for labeling based on their weighted label score. These **two sets of N and R are combined** with a user-defined ratio to yield the final set of images to label for a user in a given session.
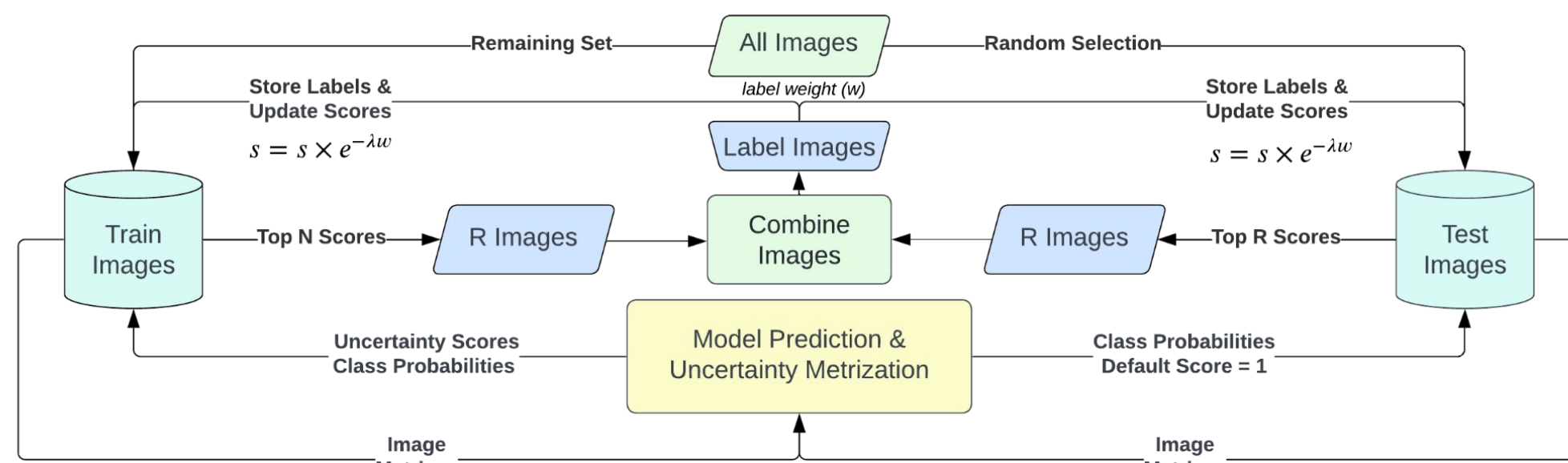


**Figure 4**: Selection of Images to be Labeled
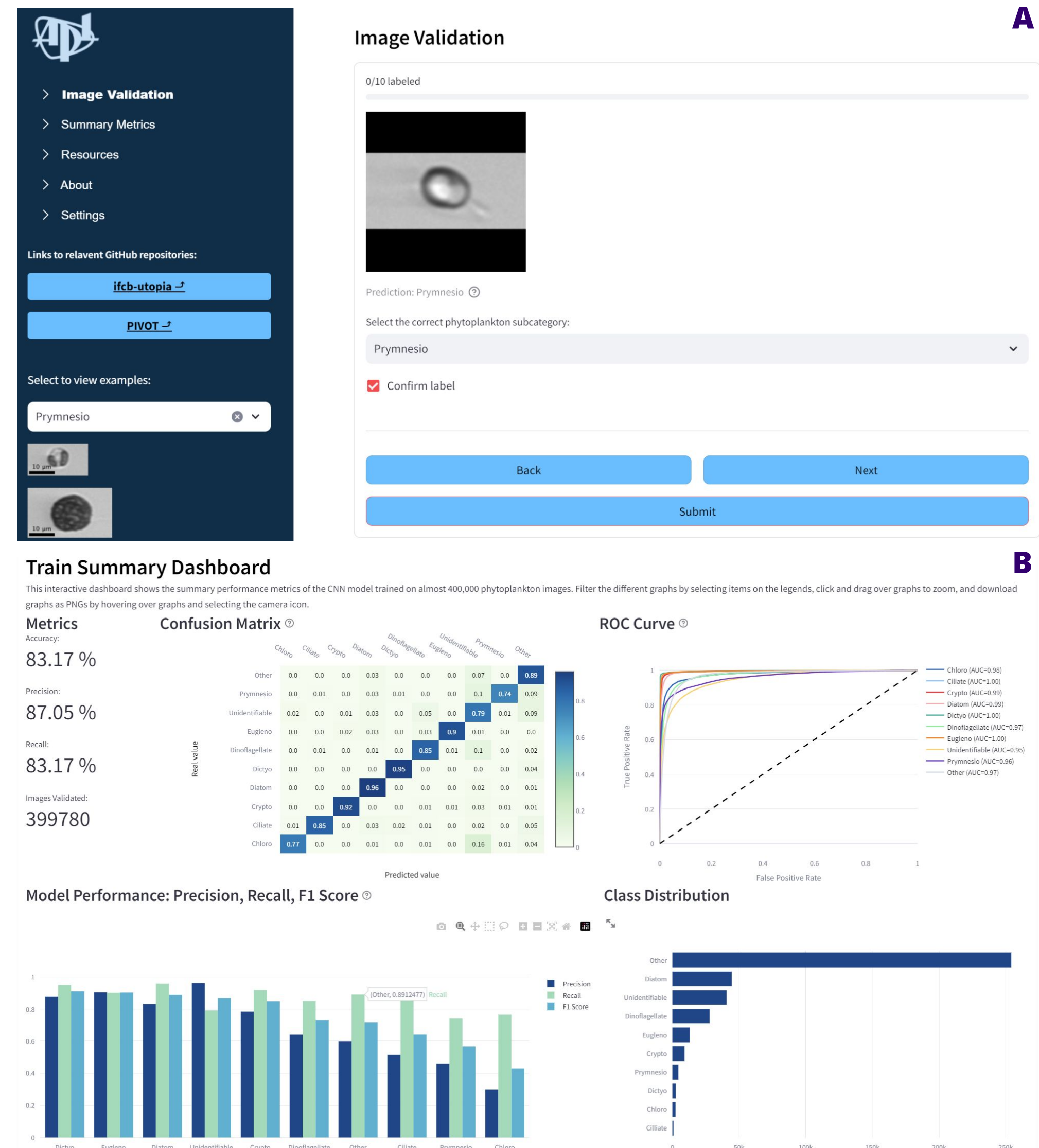
## Results



**Figure 5**: Image Validation Tool (A) and Summary Metrics Dashboard: Train Set (B)

## Conclusion

PIVOT addresses the absence of a robust platform for efficiently validating CNN models for phytoplankton image classification. It tackles the challenges of systematically processing and refining models used to label millions of single-celled autotroph images for oceanographic research.  It serves as **a central hub for researchers** to validate ML model predictions and **strategically recommends images for maximum model improvement**.

Leveraging Streamlit and Azure's underlying compute, **we centralize data storage and model processing** while **ensuring future adaptability**. Demonstrably, PIVOT paves the way for connection to new models in future research endeavors. Still, work remains:

- Investigating additional Active Learning (AL) query strategies
- Optimizing storage of intermediate data required for AL querying
- Model Fine Tuning within in the AL loop