

Plankton Image Recognition

Group 11

M102040016李翊瑄 M102040010許育禎



Content

- **Introduction**
- **Data Preprocessing and Splitting**
- **Models**
- **Insight and Future Work**

1

Introduction



Introduction -- Kaggle Competition

Featured Prediction Competition

National Data Science Bowl

Predict ocean health, one plankton at a time

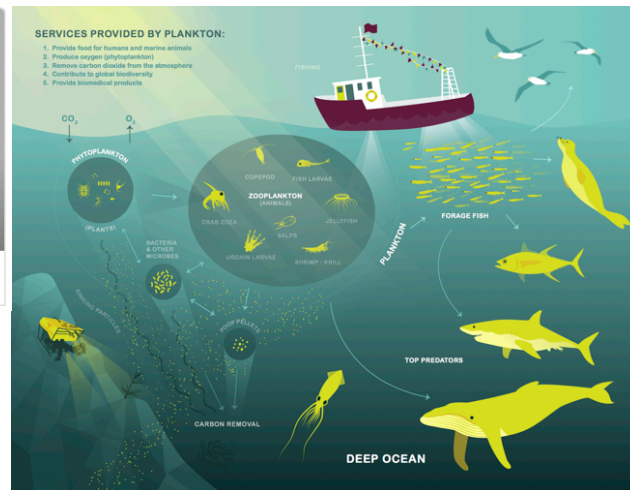
Booz Allen | Booz Allen Hamilton · 1,049 teams · 7 years ago

[Overview](#) [Data](#) [Code](#) [Discussion](#) [Leaderboard](#) [Rules](#)

[Join Competition](#)

\$175,000
Prize Money

“ Identify images of 121 species of plankton,
and calculate the probability of each species. ”





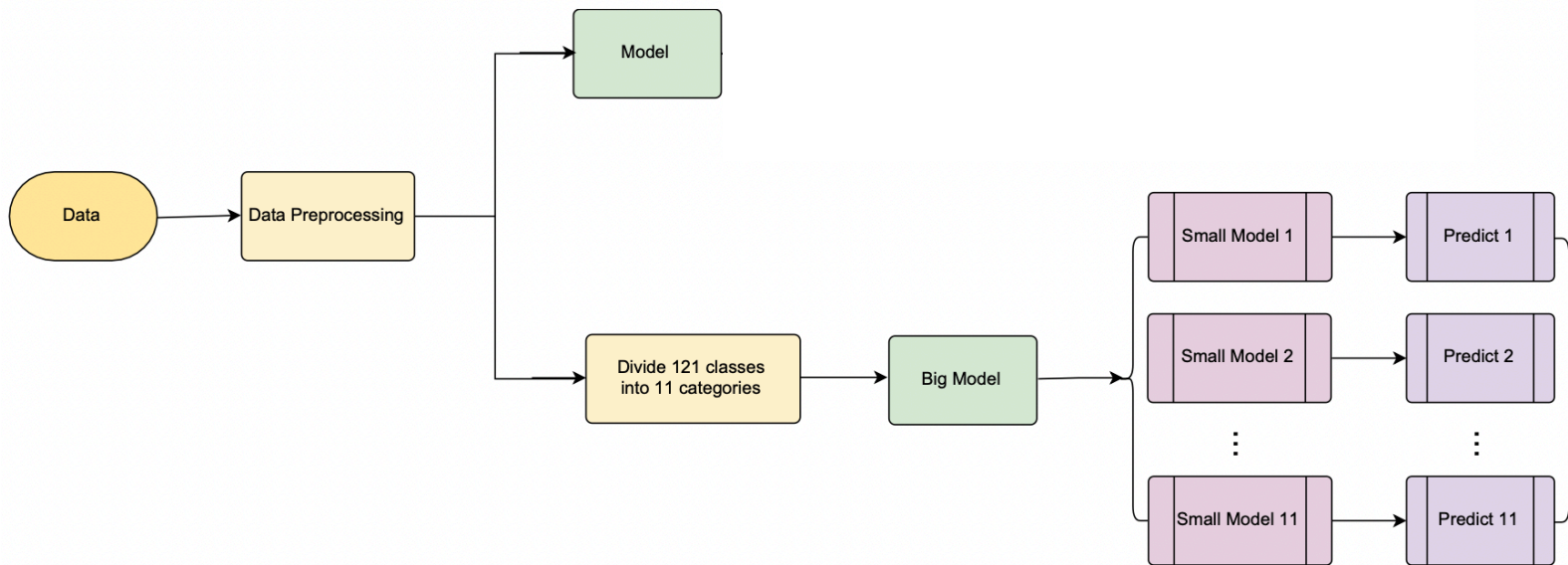
Evaluation

Image	Class 1	Class 2	Class 121
Image1.jpg	0.01	0.03	0.09
Image2.jpg	0.05	-0.01	0.02

$$\text{logloss} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log(p_{ij}),$$



Flow Chart



2

Data Preprocessing and Splitting



Data Preprocessing and Splitting

- The datasets is monochromic image -> change 1 layer to 3 layer
- Enhance the edges with ``ImageFilter.EDGE_ENHANCE_MORE`` in PIL package
- Resize to (256,256)
- Split Training set (n=24322) and Val set(n = 6014) with 80/20 principal, Test set(n=130400)



Before



After

3

Models

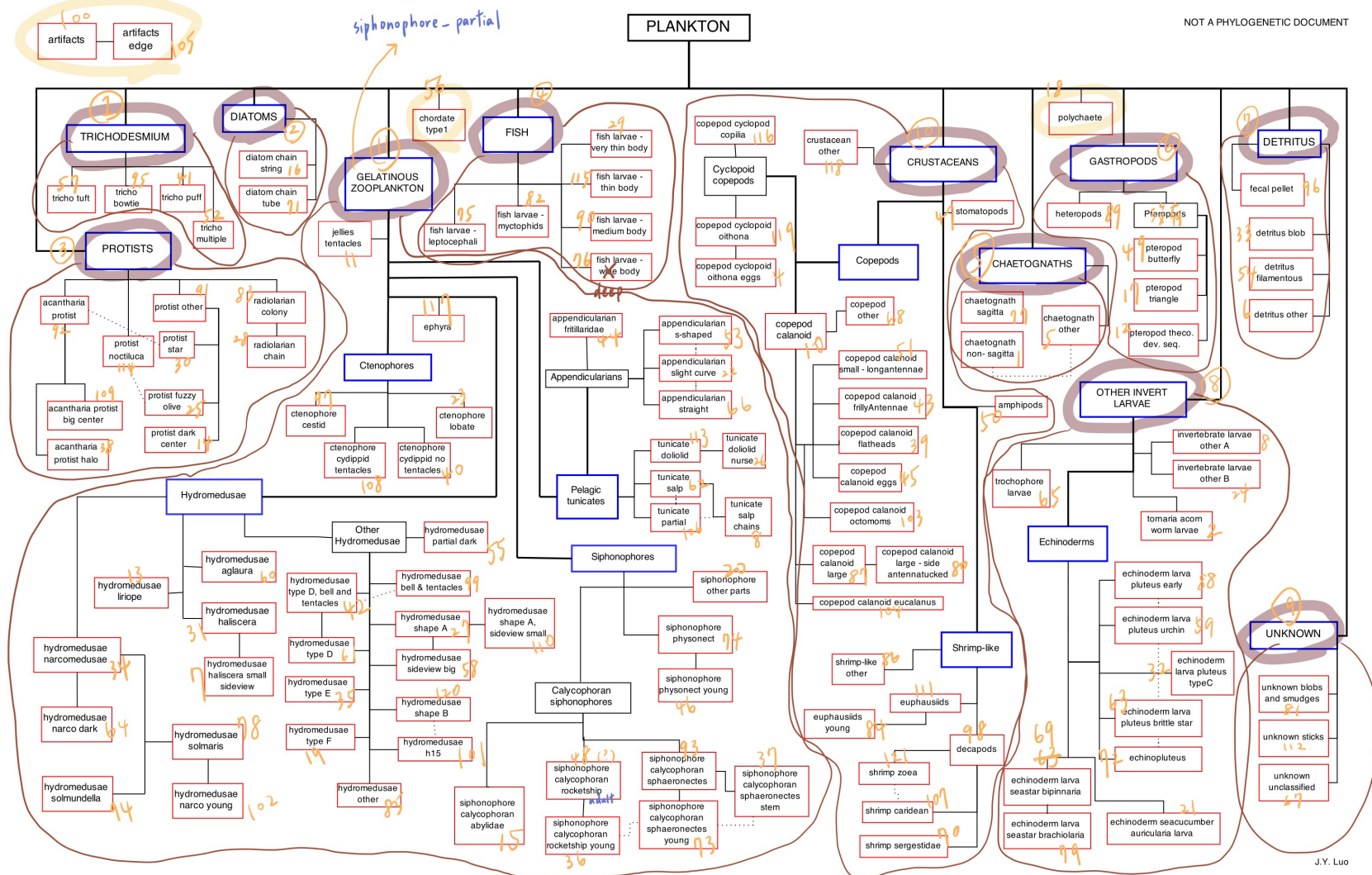


Model

Model	Optimizer	Learning Rate	Best Accuracy
CNN-1	SGD (momentum = 0.9)	10^{-2}	0.55
CNN-2	Adam (betas=(0.9, 0.999))	10^{-5}	0.65
CNN-2	SGD (momentum = 0.8)	10^{-4}	0.51
ResNet34 (pretrained = False)	SGD (momentum = 0.9)	10^{-2}	0.63
ResNet152 (pretrained = False)	SGD (momentum = 0.9)	10^{-2}	0.68



Kaggle Score: 10.88





Big Model -> Small Model

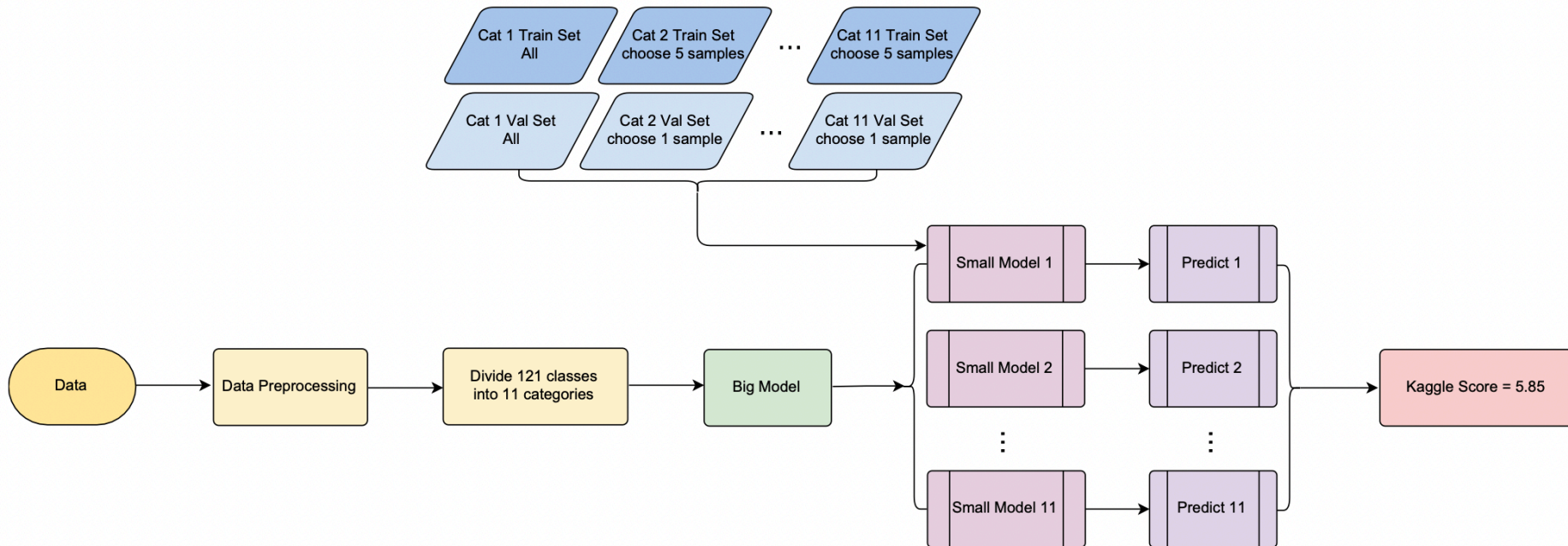
Big Model

Split all Classes to 11 Categories in train set and validation set.

file_path	images_category	label	label_11	label2
protist_star_90.jpg	protist_star	30	protists	3
protist_star_78.jpg	protist_star	30	protists	3
:	:	:	:	:
artifacts_399.jpg	artifacts	100	unknown	9
artifacts_197.jpg	artifacts	100	unknown	9
:	:	:	:	:
hydromedusae_other_117.jpg	hydromedusae_other	85	gelatinous_zooplankton	11
hydromedusae_other_115.jpg	hydromedusae_other	85	gelatinous_zooplankton	11
:	:	:	:	:
copepod_calanoid_1097.jpg	copepod_calanoid	10	crustaceans	10
copepod_calanoid_639.jpg	copepod_calanoid	10	crustaceans	10



Big Model -> Small Model -- Flow Chart





Big Model -> Small Model

Take 1st small model for example:

Train Set

Get all training images in 1st Category, and randomly pick up 5 images in the rest of classes

Validation Set

Get all validation images in 1st Category, and randomly pick up 1 image in the rest of classes

Test Set

All test images predicted as 1st Category



Comparison

Model	Kaggle Score
Model (ResNet152)	10.88
Big Model -> Small Model (ResNet152) -> (ResNet152)	5.85



Depends on the accuracy of big model



New Findings in Small Model Training

- Val Acc < Training Acc in the beginning of the training
- Sometimes accuracy improves, but log loss does not

```
===== Epoch 1 =====  
Train Acc: 0.452481 Train Loss: 2.765803  
Val Acc: 0.565217 Val Loss: 1.276544  
===== Epoch 2 =====  
Train Acc: 0.524806 Train Loss: 2.040130  
Val Acc: 0.668116 Val Loss: 1.064313  
===== Epoch 3 =====  
Train Acc: 0.528093 Train Loss: 1.902031  
Val Acc: 0.475362 Val Loss: 1.413311  
===== Epoch 4 =====  
Train Acc: 0.525403 Train Loss: 1.815078  
Val Acc: 0.626087 Val Loss: 0.998090
```

```
===== Epoch 20 =====  
Train Acc: 0.633592 Train Loss: 1.212397  
Val Acc: 0.614493 Val Loss: 1.036882  
===== Epoch 21 =====  
Train Acc: 0.637179 Train Loss: 1.223277  
Val Acc: 0.668116 Val Loss: 1.235346  
===== Epoch 24 =====  
Train Acc: 0.663479 Train Loss: 1.126488  
Val Acc: 0.591304 Val Loss: 1.582351  
===== Epoch 25 =====  
Train Acc: 0.661088 Train Loss: 1.119432  
Val Acc: 0.550725 Val Loss: 1.416806
```


5

Insight and Future Work



Insight and Future Work



Whether plankton unrelated to other species be included in the category of unknown species

2

Apply Semi-supervised Learning method.

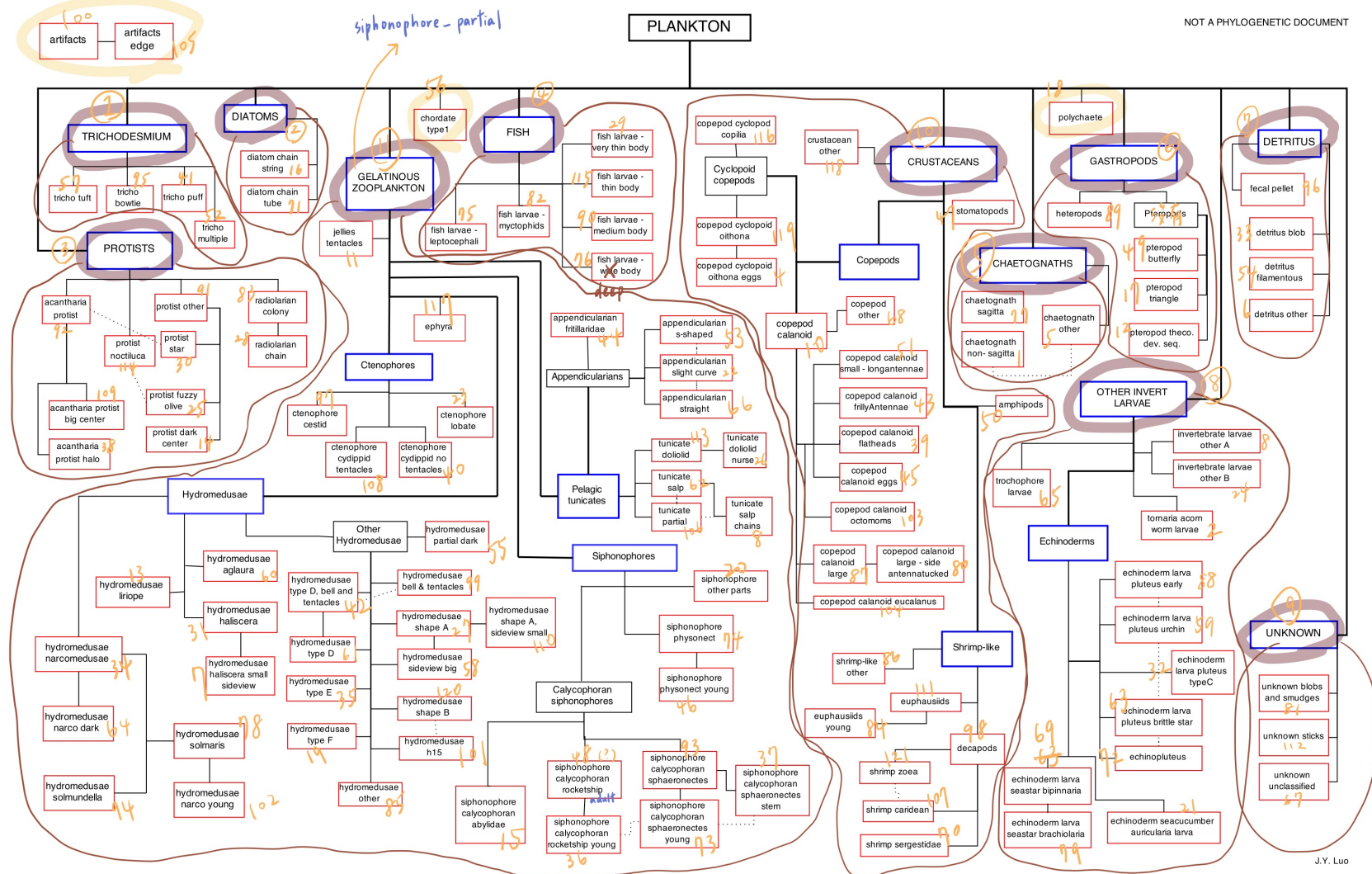
4

1

After being divided into 11 categories, the amount of data in each category is unbalanced.
Categories with too many species can be subdivided.

3

The way to form a new training set and validation set.





THANKS!