

# Applied ML Final Project

## Detection of Insolubles in Real Ice Core Data

Amalie Mygind, Marcus Nygaard,  
Michala Jensen, Søren Langkilde  
& Ulrik Friis-Jensen

UNIVERSITY OF COPENHAGEN



# Outline

- Introduction
- Presentation of models
  - K-means
  - NN
  - CNN
  - Ensemble model
  - Mixed model
- Evaluation of models
- Prediction on real data



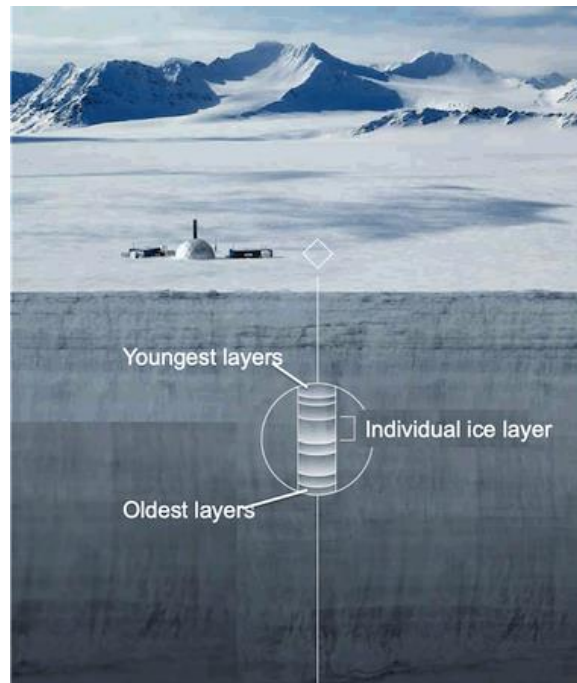
Image:

[http://lindseynicholson.org/wp-content/uploads/2018/04/1\\_Core\\_mount-hunter-ice-core-preserved-climate-and-environmental-record-eos\\_org.jpg](http://lindseynicholson.org/wp-content/uploads/2018/04/1_Core_mount-hunter-ice-core-preserved-climate-and-environmental-record-eos_org.jpg)

# Introduction

- Goal: To classify particles in ice core data using machine learning.
- Approach: Utilizing all available data to create machine learning algorithms for classification and combining them in efforts to improve accuracy.

# The ice core data



1 datapoint  
=  
1 image + meta data

Image:

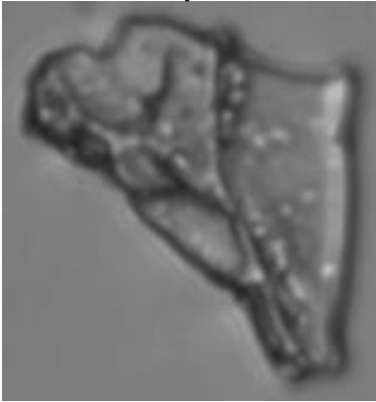
<https://theconversation.com/explainer-what-are-ice-cores-24302>

<https://www.chemedx.org/article/ice-cores-stable-isotopes-climate-change-and-chemistry>

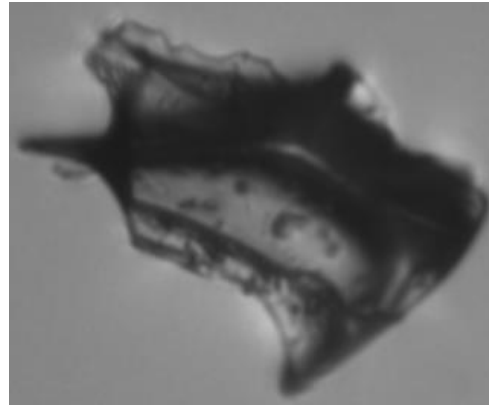
<https://www.crowcanyon.org/index.php/ice-core-studies>

# The data used for training

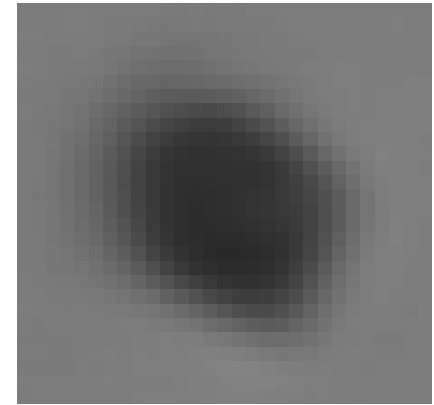
Campanian



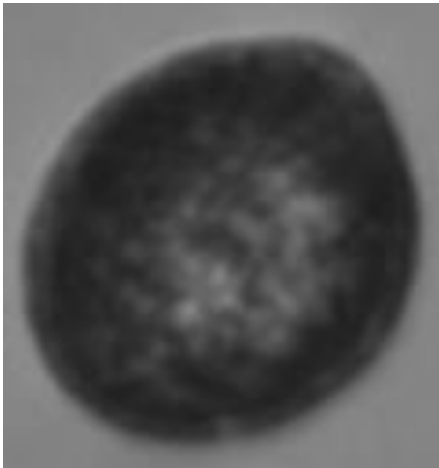
Grimsvotn



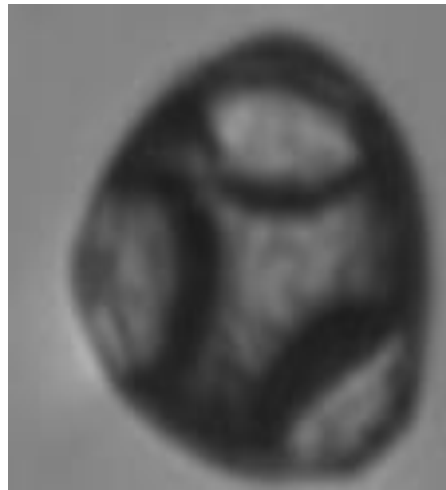
Dust



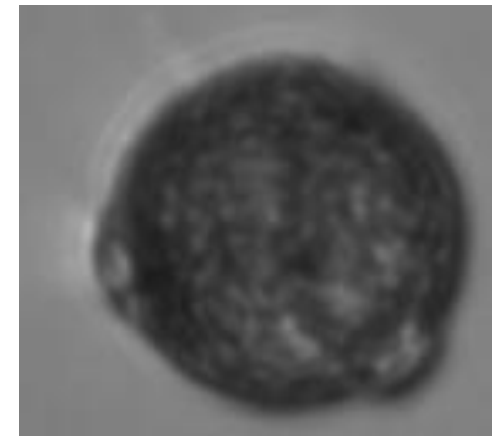
Qsuber



Corylus



Qrobur



# Meta data

	#
<b>Total features</b>	63
<b>Useable features</b>	39
<b>Unused feature</b>	24



Microscope used for data capture

# Models

K-means, NN, CNN, ensemble model & mixed model

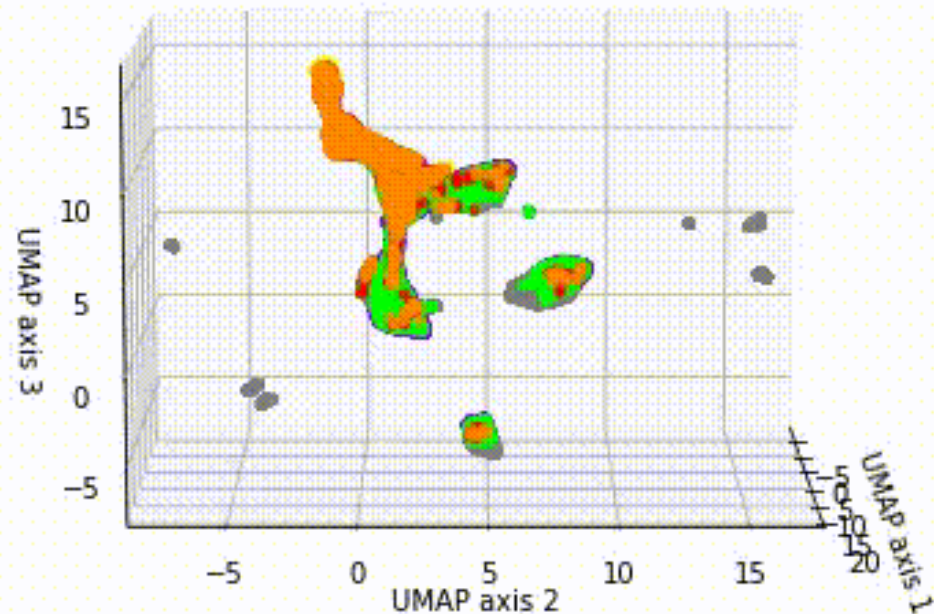
UNIVERSITY OF COPENHAGEN



# UMAP and clustering on meta data

- UMAP
  - Dimensionality reduction
  - 39 feature to 3D
  - Considered 200 neighbors
- Clustering by k-means
- Performance
  - ~ 63%

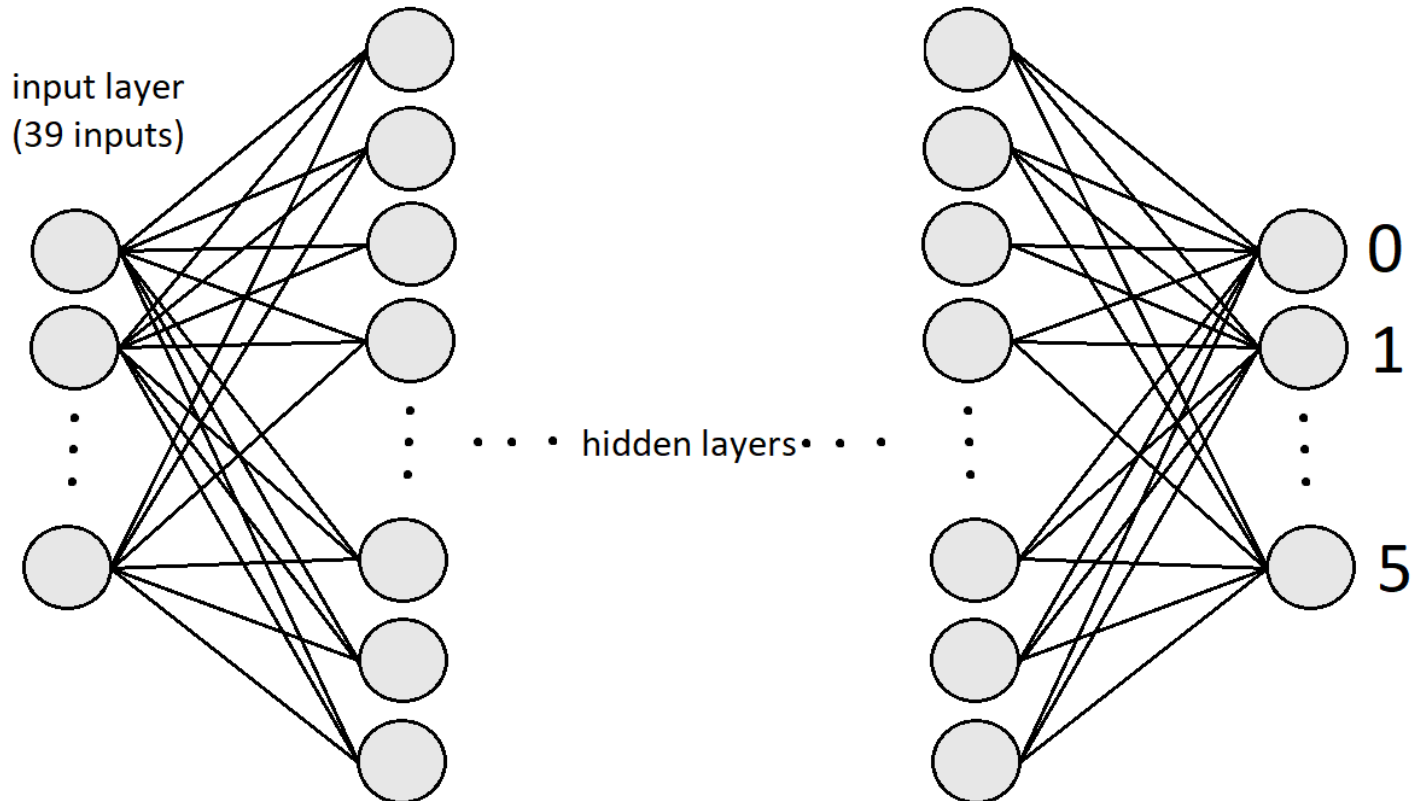
- = campanian
- = corylus
- = dust
- = grimsvotn
- = qrobur
- = qsuber





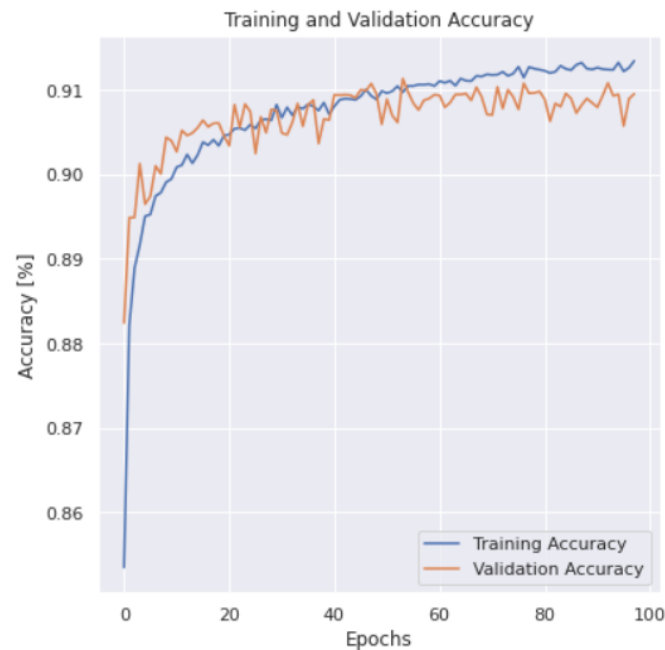
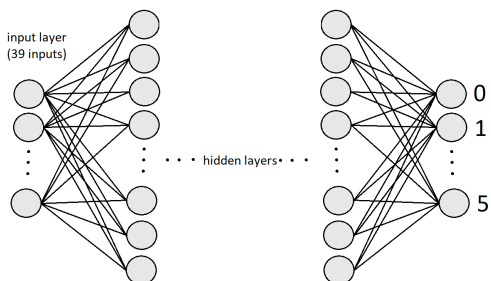
# Neural Network on meta data

- 39 inputs, 4 hidden layers, 50 neurons in each



# Neural Network on meta data

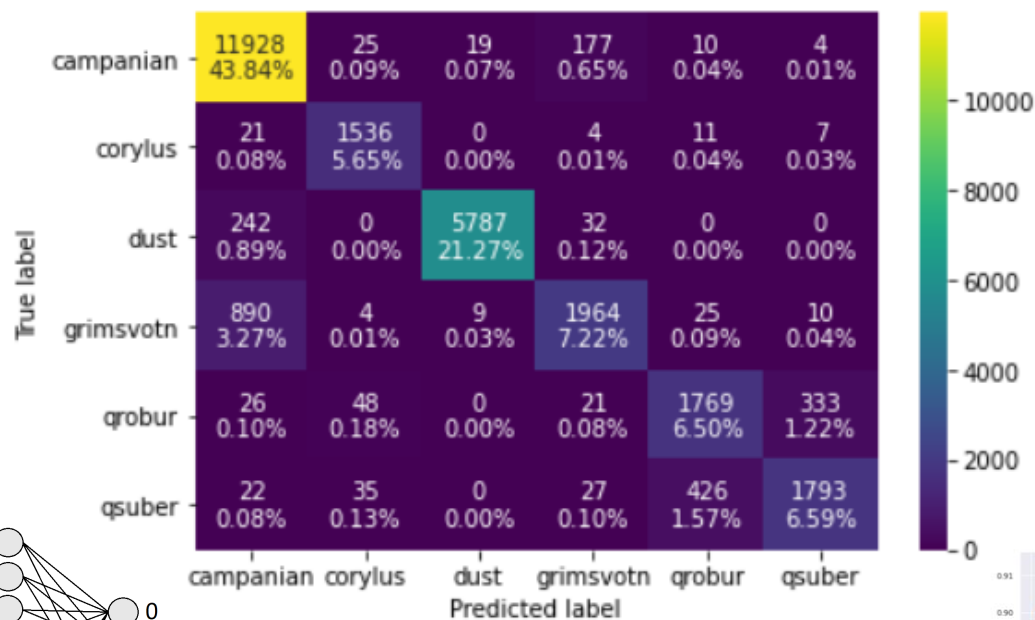
- 39 inputs, 4 hidden layers, 50 neurons in each
- Performance
  - Almost perfect dust predictor
- Loss function: 0.26 (categorical cross entropy)



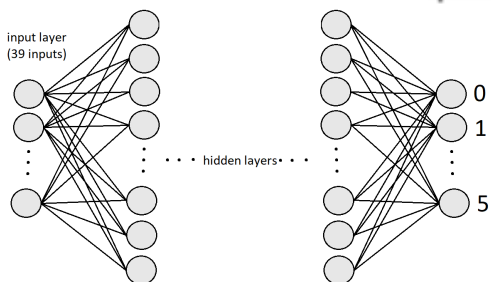
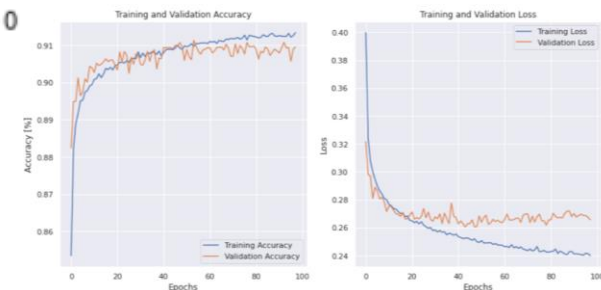
# Neural Network on meta data

- 39 inputs, 4 hidden layers, 50 neurons in each
- Performance
  - Almost perfect dust predictor
- Loss function: 0.26 (categorical cross entropy)

NN on meta data

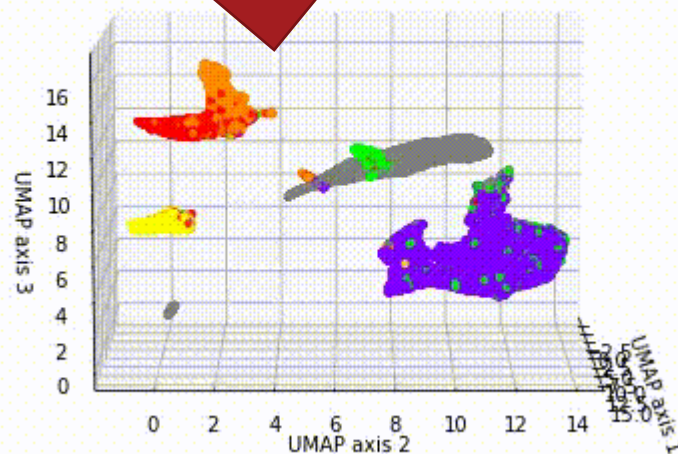
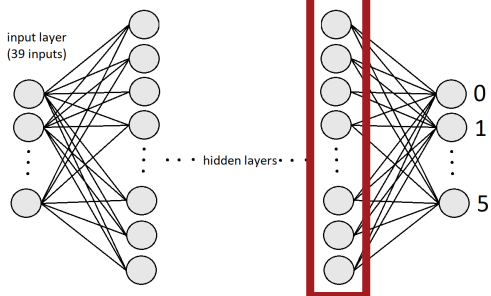
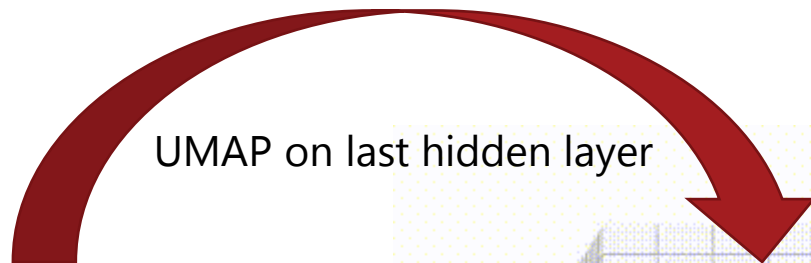


Accuracy=91.08%



# Neural Network on meta data

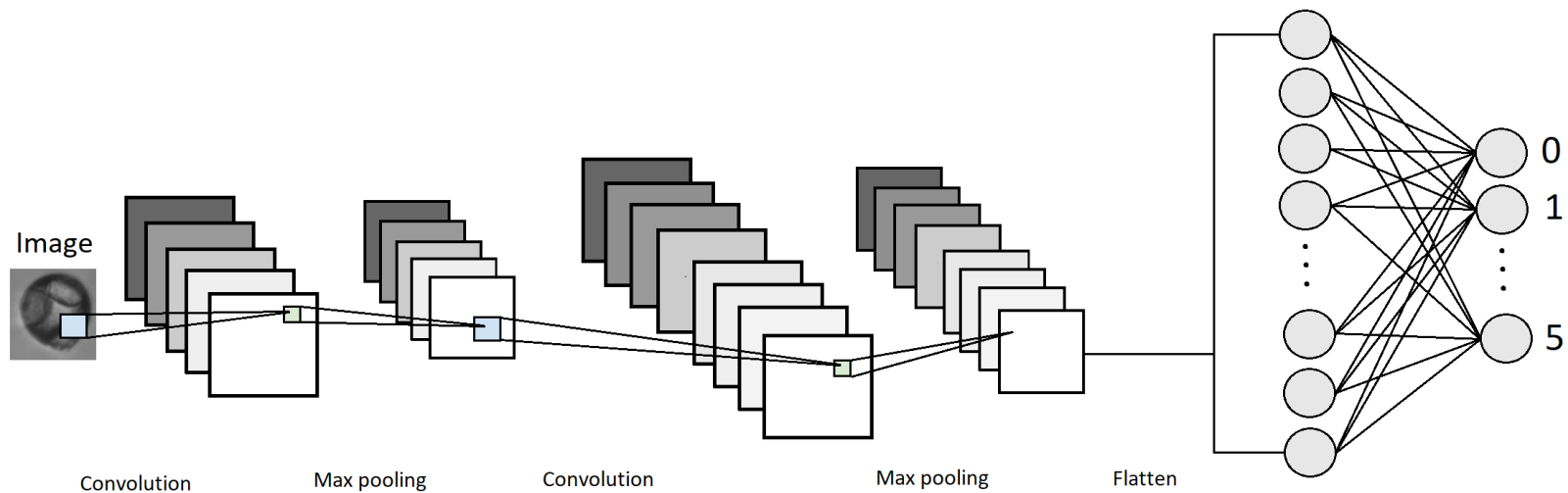
- 39 inputs, 4 hidden layers, 50 neurons in each
- Performance
  - Almost perfect dust predictor
- Loss function: 0.26 (categorical cross entropy)
- Primary classes group together



- = campanian
- = corylus
- = dust
- = grimsvotn
- = qrobur
- = qsuber

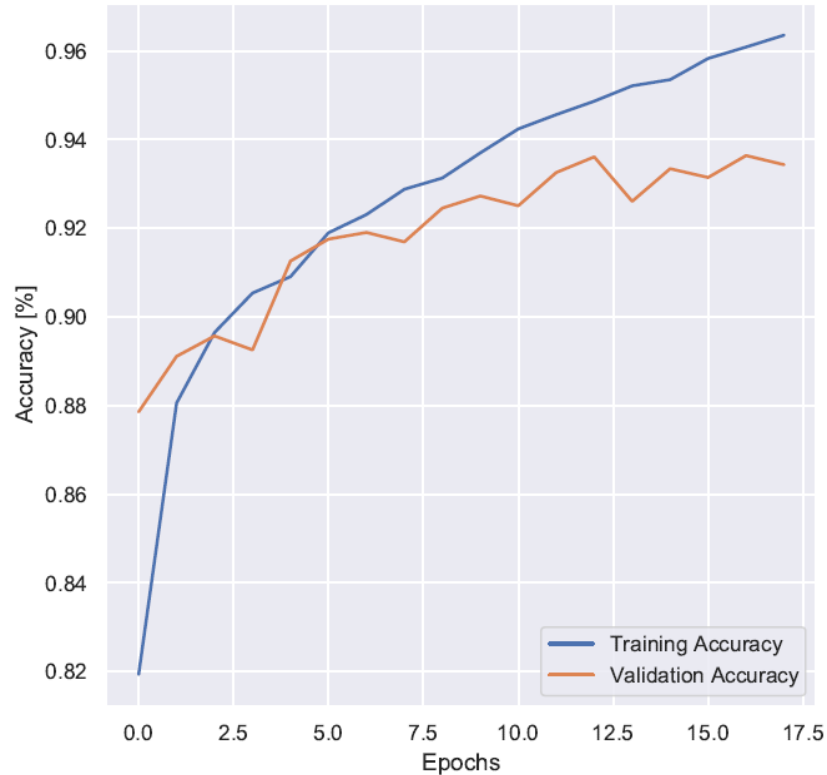
# CNN on image data

- Convolutional Neural Network
  - 4 conv layers w/ max pool, 2 drop out layers
- Bayesian optimization of hyperparameters
  - # filters, drop out rate, kernel regularizer factor

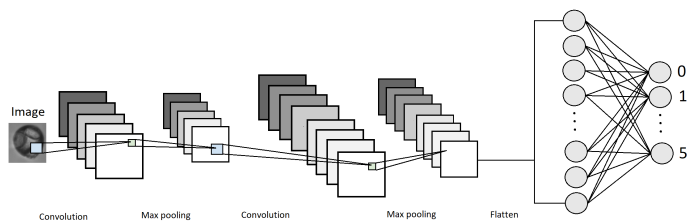
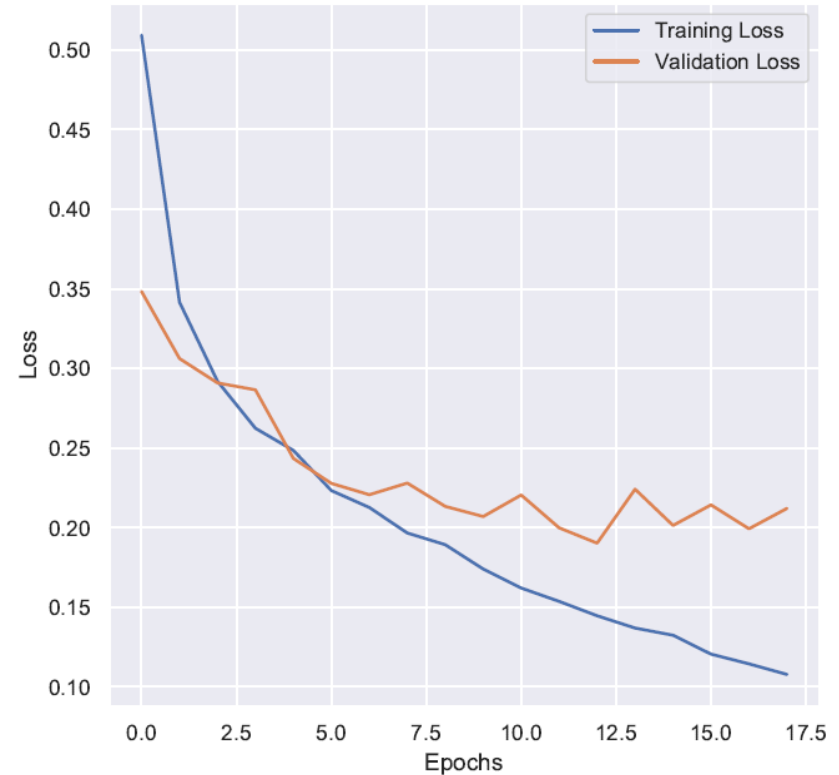


# CNN on image data

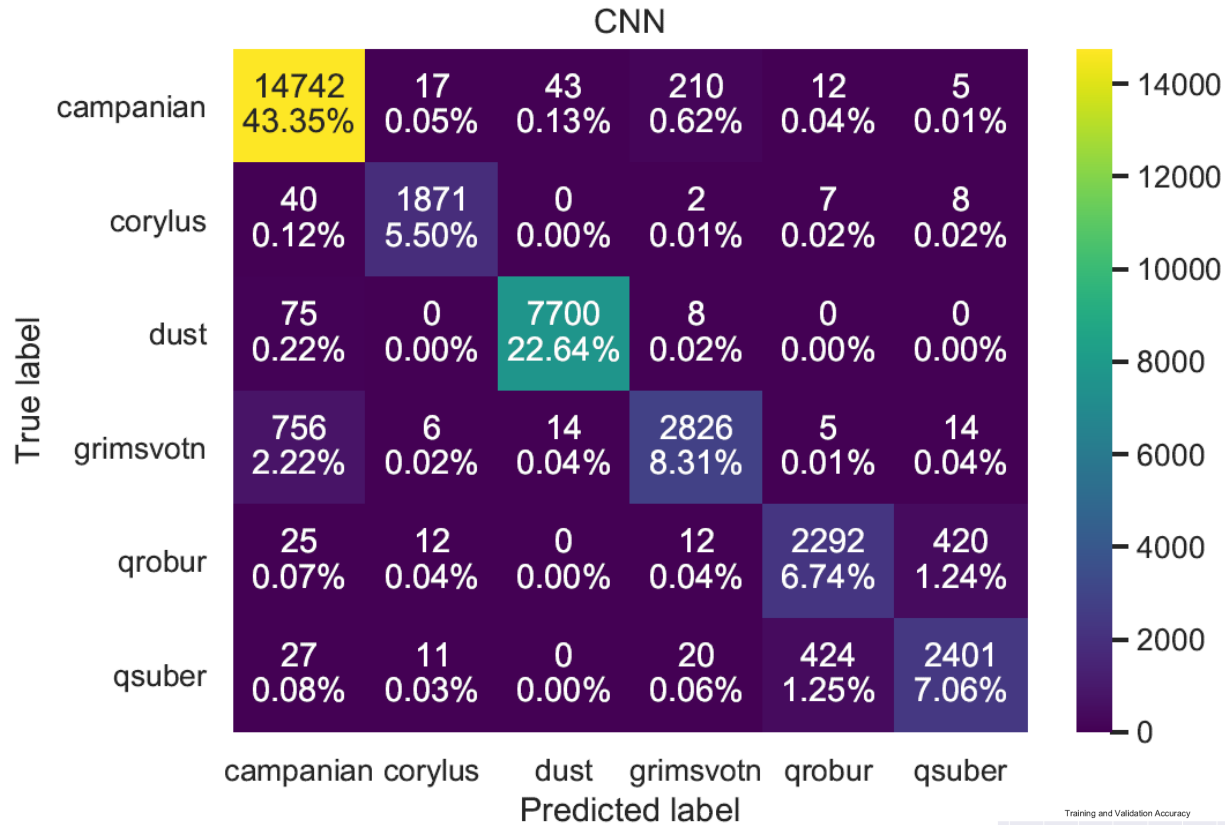
### Training and Validation Accuracy



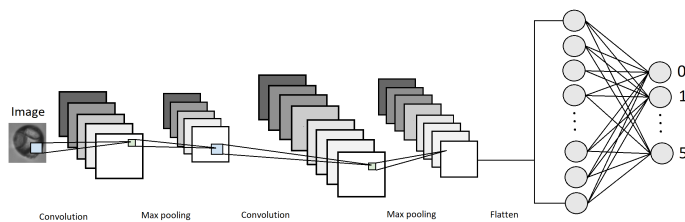
### Training and Validation Loss



# CNN on image data

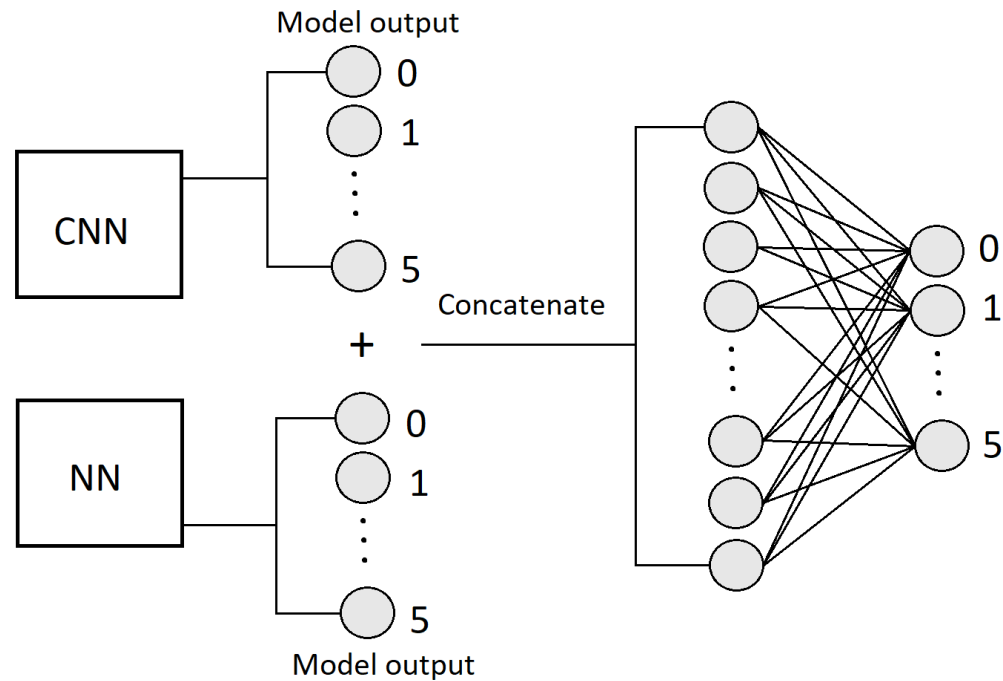


Accuracy=93.61%



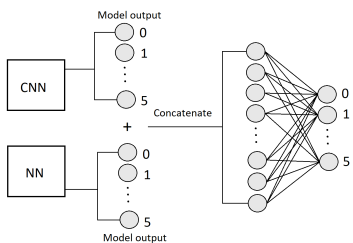
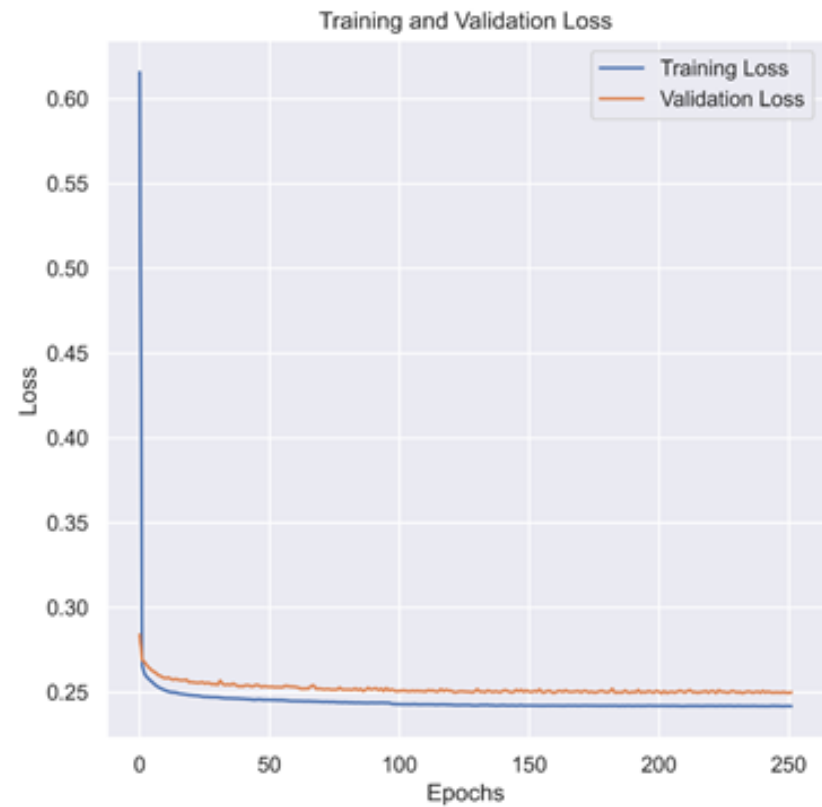
# Ensemble model

- Combine output from CNN and NN into a single input vector

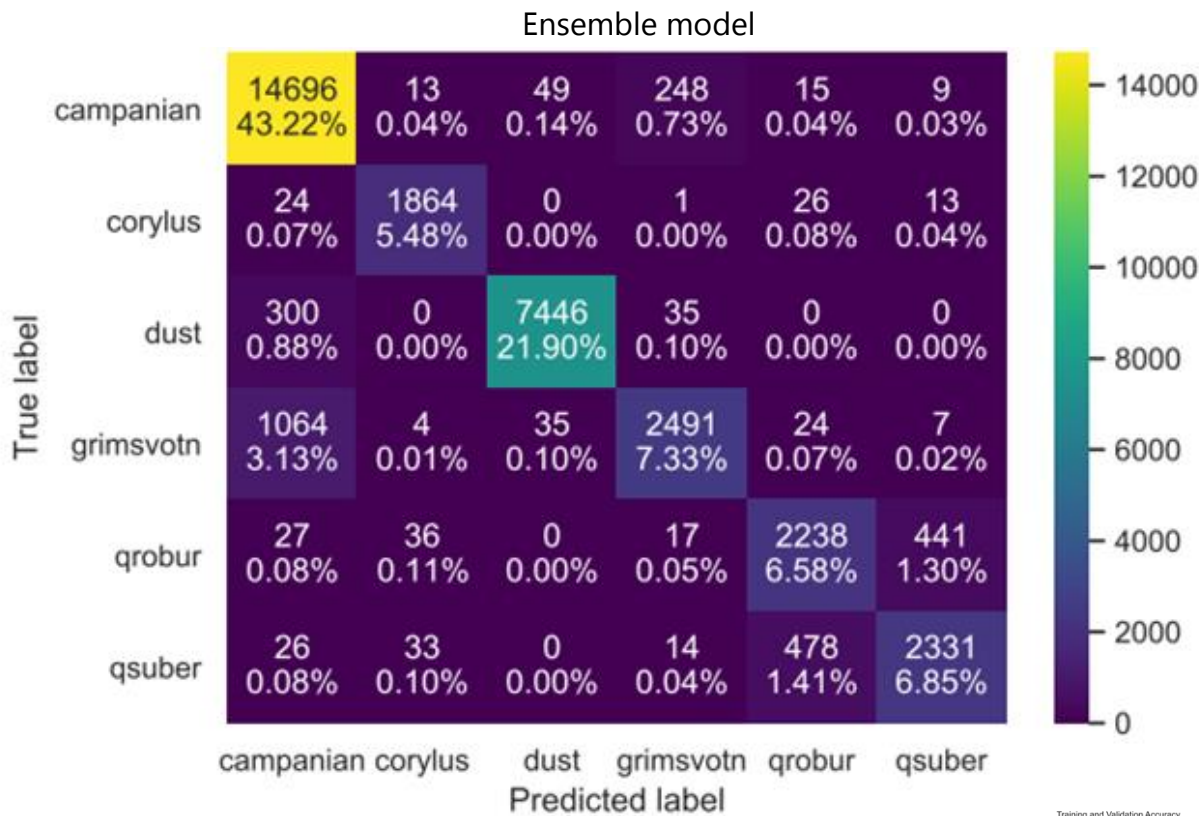




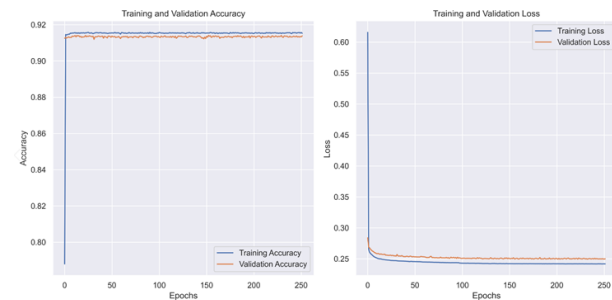
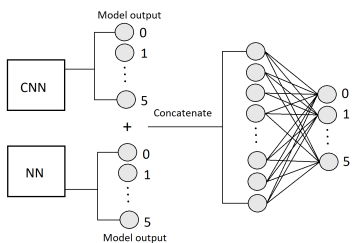
# Ensemble model



# Ensemble model

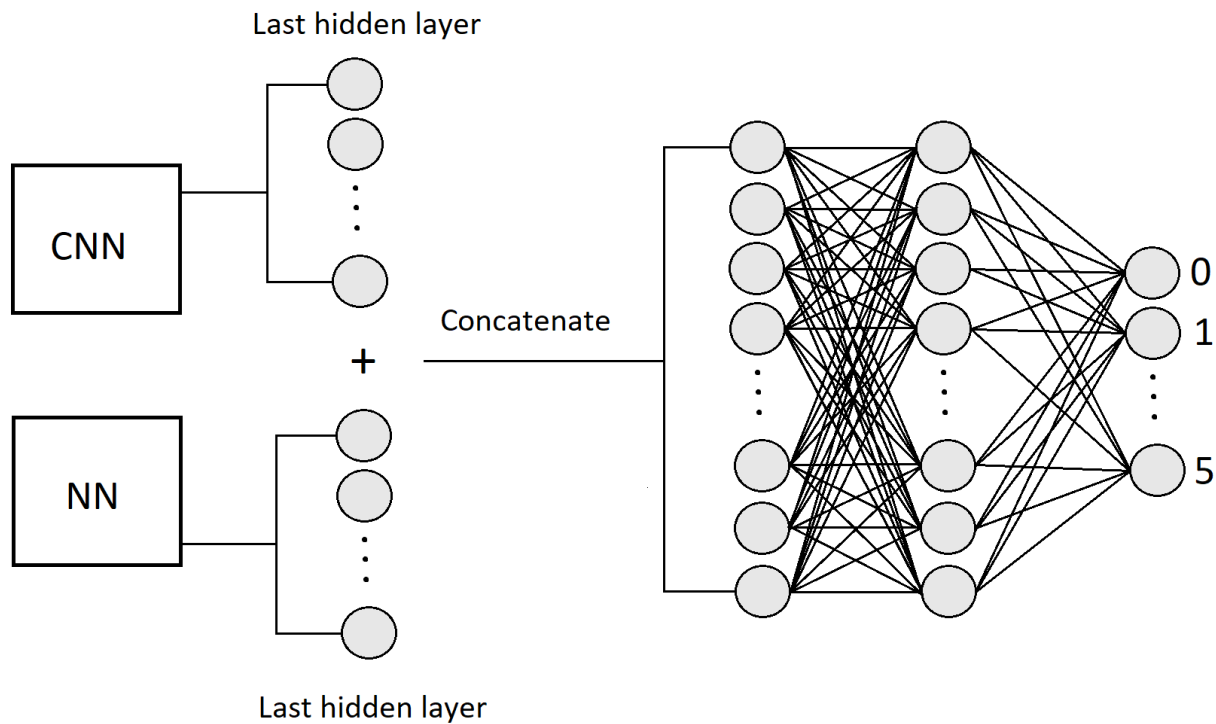


Accuracy=91.36%



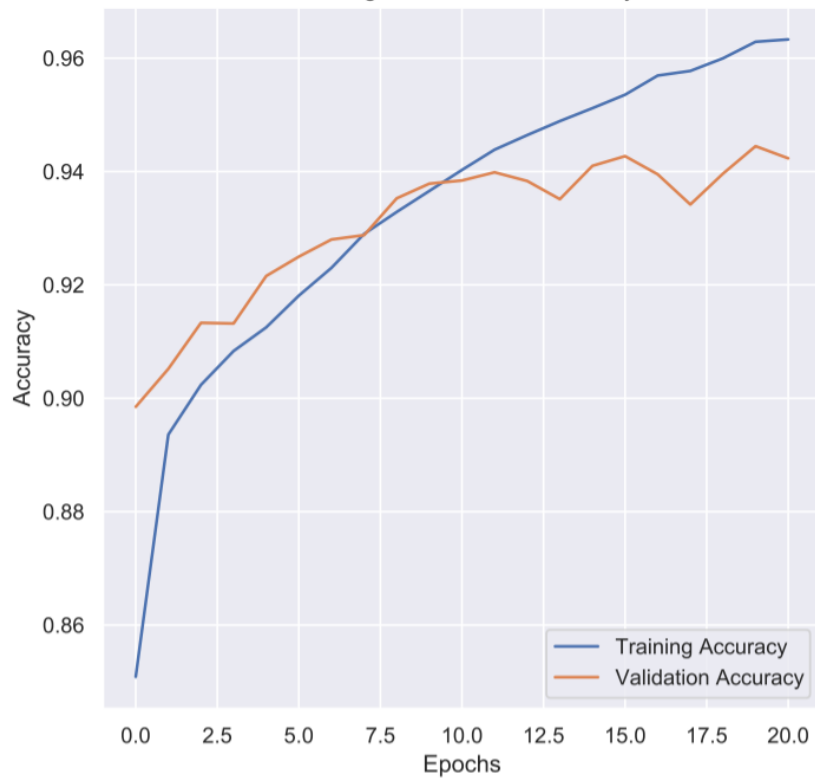
# Mixed model

- Two inputs
  - Image data
  - Meta data
- Custom data loader to handle different data types

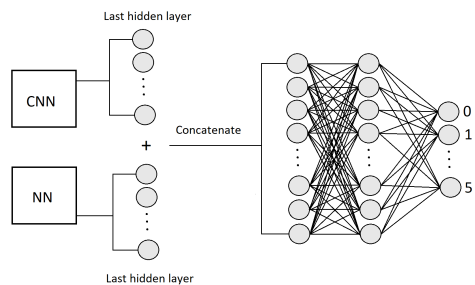
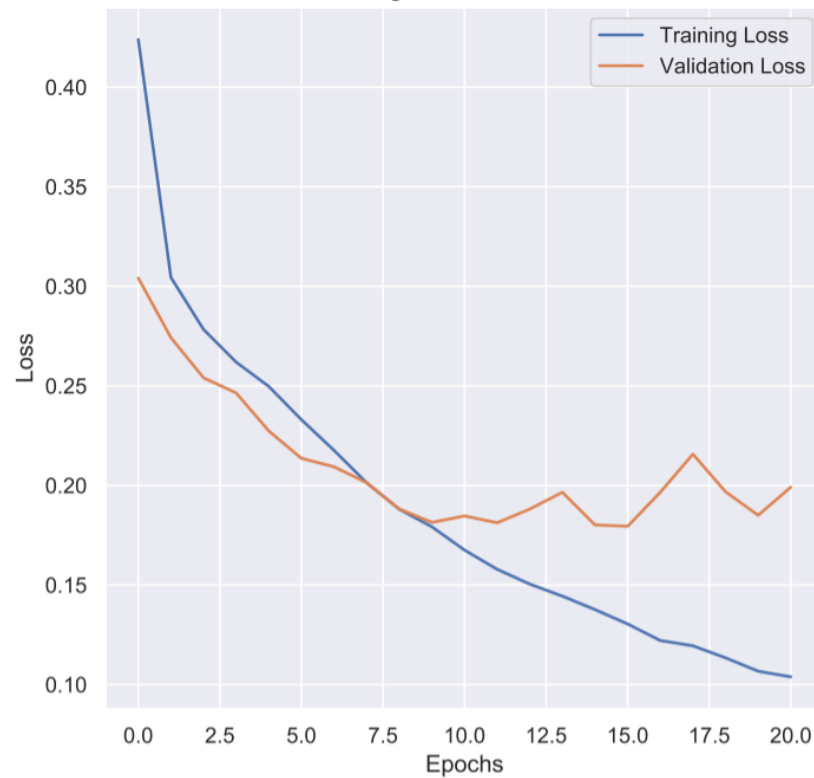


# Mixed model

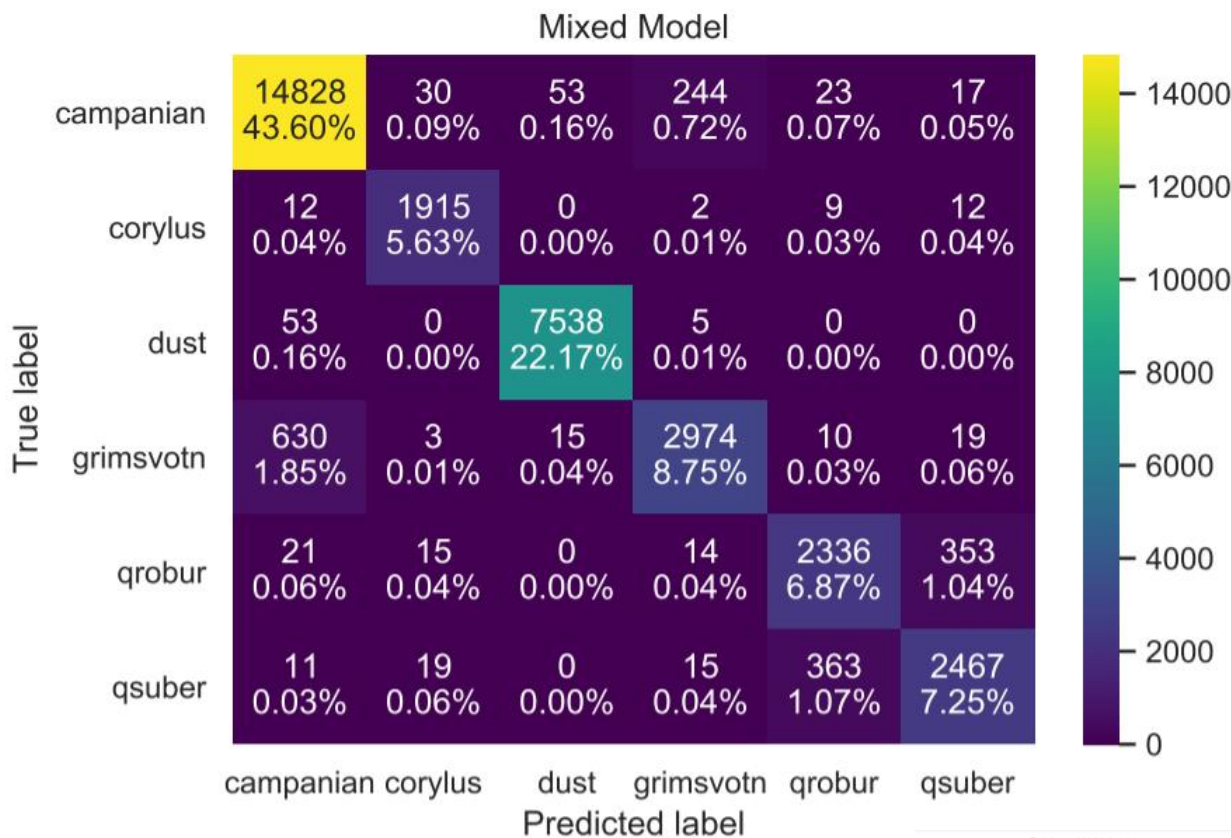
Training and Validation Accuracy



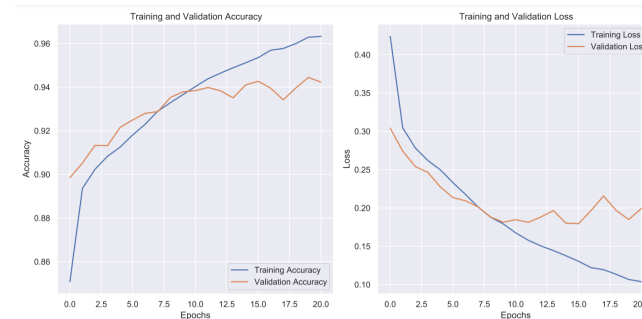
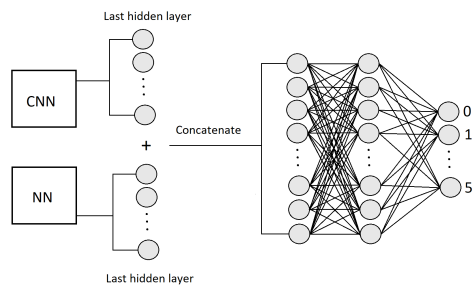
Training and Validation Loss



# Mixed model



Accuracy=94.27%



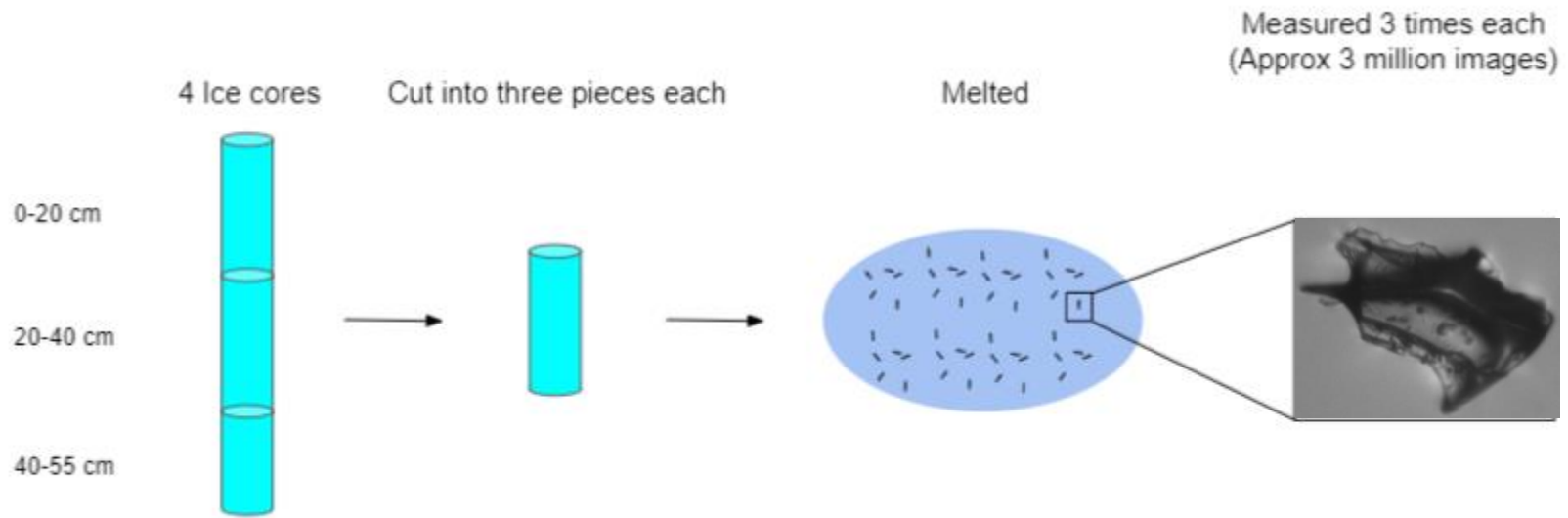
Type of model	Trainable parameters	Cross Entropy Loss	Accuracy	Evaluation
<b>NN</b>	9,956	0.26	91.08%	Simple, Medium performance
<b>CNN</b>	4,022,554	0.19	93.61%	High complexity, Good performance
<b>Mixed model</b>	4,061,098	0.18	94.27%	High complexity, Best performance
<b>Ensemble model</b>	276	0.25	91.36%	Simple (Requires other models), Medium performance

# Real ice core data

UNIVERSITY OF COPENHAGEN



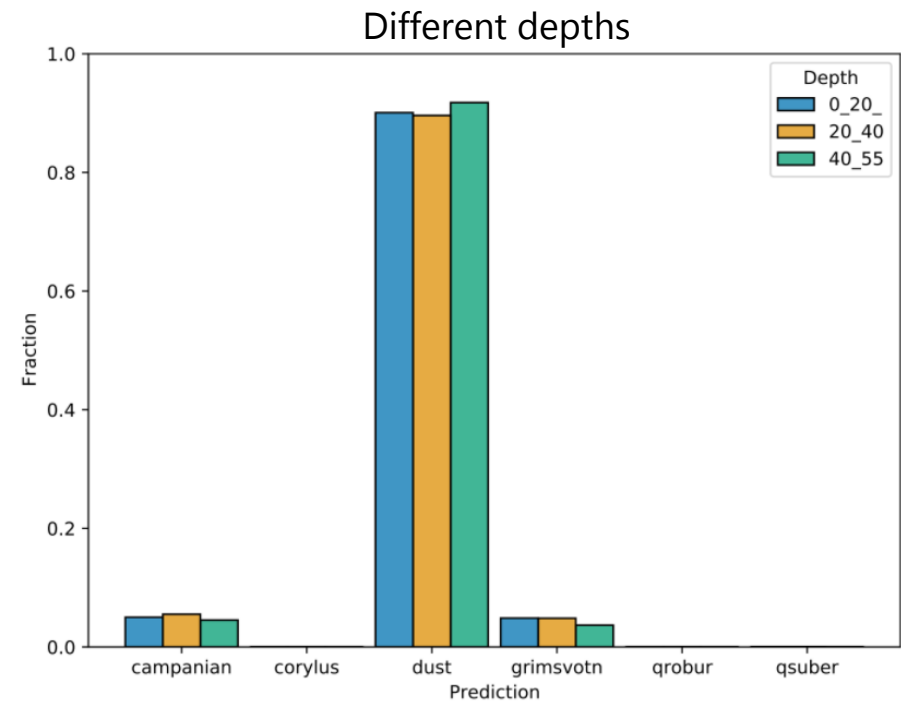
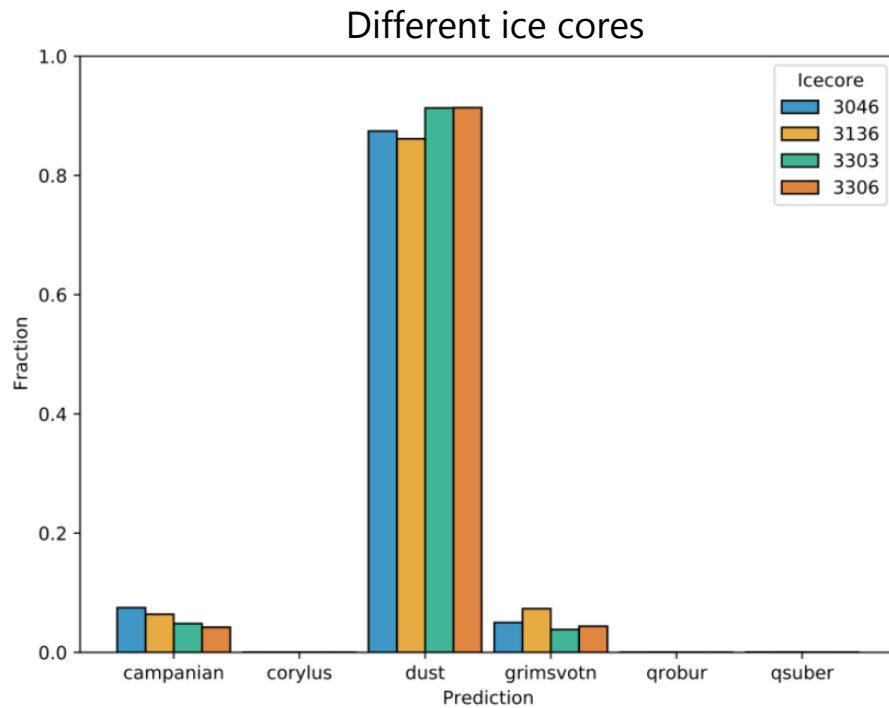
# The ice core data used for testing



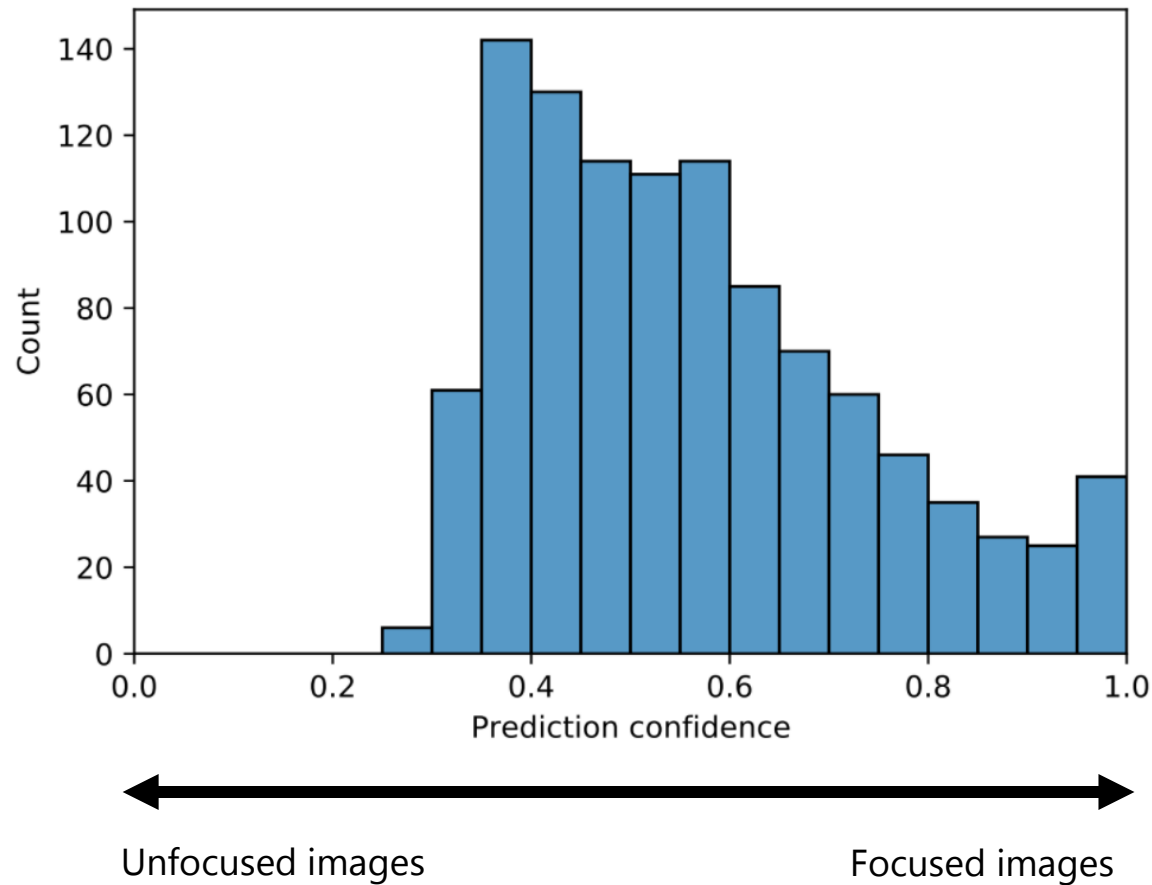


# Model predictions on test data

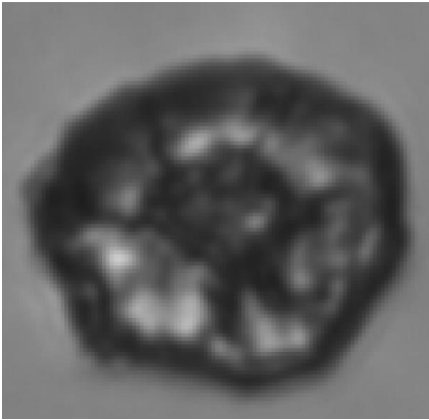
- Predictions using mixed model
- Expected 99% dust, 1% ash/other, 0% pollen
- Predicted 90.4% dust, 9.6% ash, 0.035% pollen



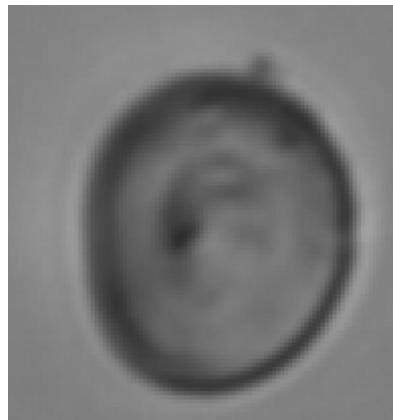
# Confidence level on pollen predictions



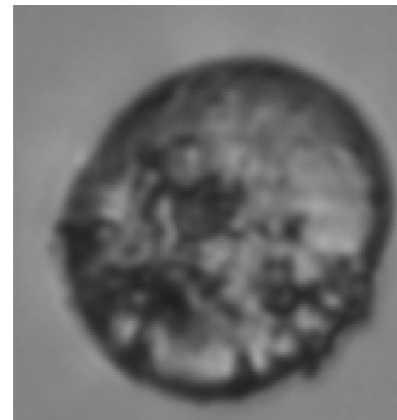
# 6 most interesting images in test data



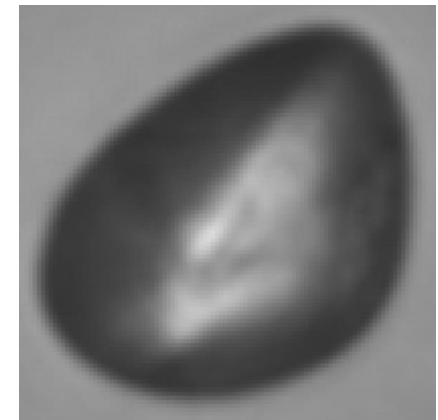
test\_GRIP\_31may2021/GRIP\_raw/GRIP\_3136\_20\_40\_3/GRIP\_3136\_20\_40\_3\_16380.png



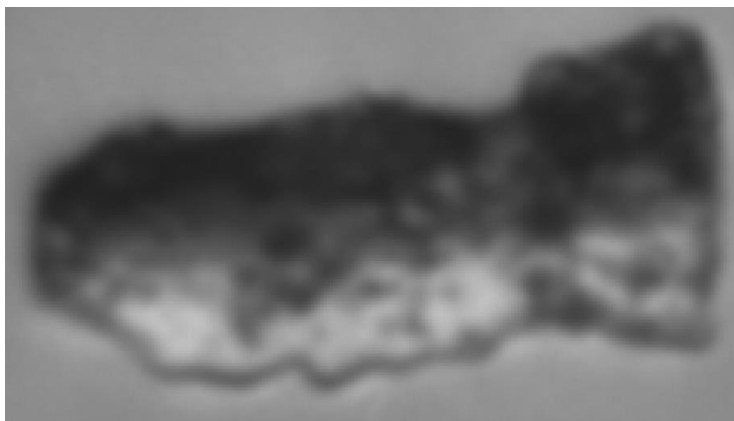
test\_GRIP\_31may2021/GRIP\_raw/GRIP\_3306\_0\_20\_3/GRIP\_3306\_0\_20\_3\_100001.png



test\_GRIP\_31may2021/GRIP\_raw/GRIP\_3136\_40\_55\_3/GRIP\_3136\_40\_55\_3\_912.png



test\_GRIP\_31may2021/GRIP\_raw/GRIP\_3046\_0\_20\_1/GRIP\_3046\_0\_20\_1\_1358.png



test\_GRIP\_31may2021/GRIP\_raw/GRIP\_3136\_20\_40\_2/GRIP\_3136\_20\_40\_2\_856.png



test\_GRIP\_31may2021/GRIP\_raw/GRIP\_3136\_40\_55\_2/GRIP\_3136\_40\_55\_2\_3767.png

# Thank you for listening!

## Any questions?



# Appendix

UNIVERSITY OF COPENHAGEN

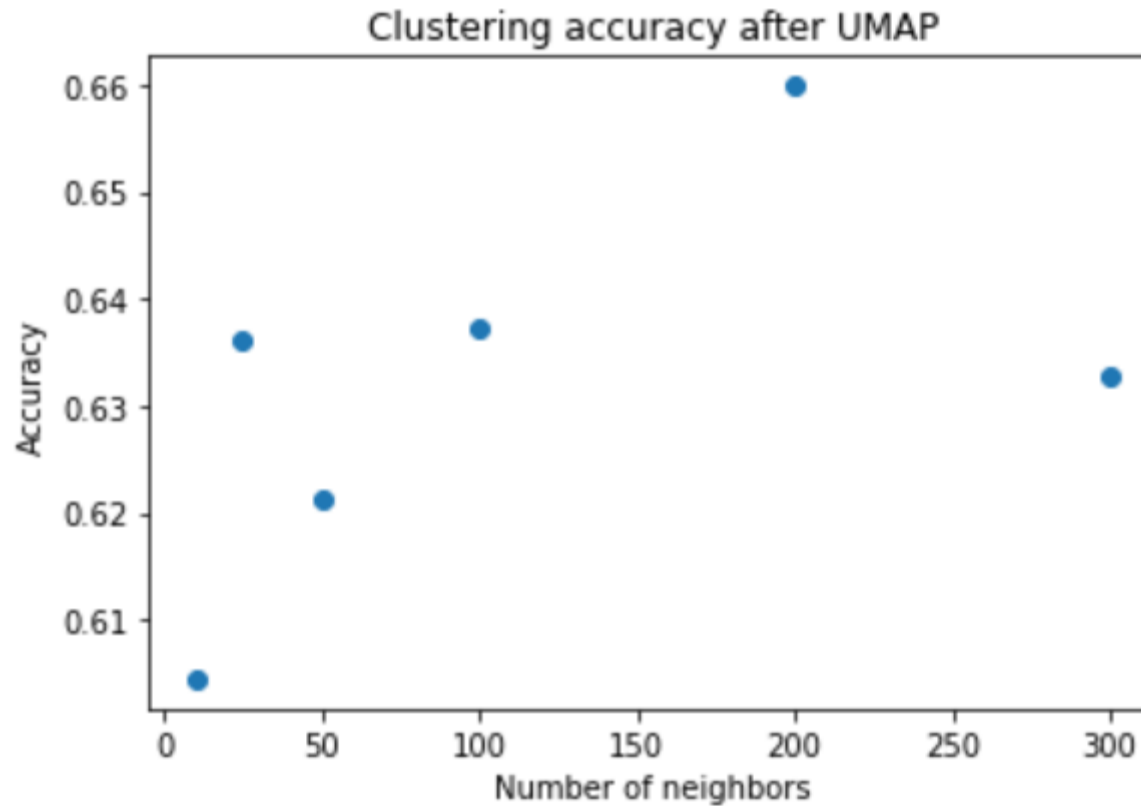


# Which meta data features were dropped/ used

```
dropped_features=['Particle ID', 'Calibration Image', 'Camera', 'Capture X',  
'Capture Y', 'Date', 'ImageFile', 'Timestamp', 'imgpaths', 'camp', 'corylus',  
'dust', 'grim', 'qrob', 'qsub', 'Time', 'Source Image', 'Image Y', 'Image X',  
'Image Height', 'Image Width', 'Filter Score', 'Elapsed Time', 'Calibration  
Factor']
```

```
used_features=['Area (ABD)', 'Area (Filled)', 'Aspect Ratio', 'Biovolume  
(Cylinder)', 'Biovolume (P. Spheroid)', 'Biovolume (Sphere)', 'Circle Fit',  
'Circularity', 'Circularity (Hu)', 'Compactness', 'Convex Perimeter', 'Convexity',  
'Diameter (ABD)', 'Diameter (ESD)', 'Edge Gradient', 'Elongation', 'Feret Angle  
Max', 'Feret Angle Min', 'Fiber Curl', 'Fiber Straightness', 'Geodesic Aspect  
Ratio', 'Geodesic Length', 'Geodesic Thickness', 'Intensity', 'Length', 'Particles  
Per Chain', 'Perimeter', 'Roughness', 'Sigma Intensity', 'Sphere Complement',  
'Sphere Count', 'Sphere Unknown', 'Sphere Volume', 'Sum Intensity', 'Symmetry',  
'Transparency', 'Volume (ABD)', 'Volume (ESD)', 'Width']
```

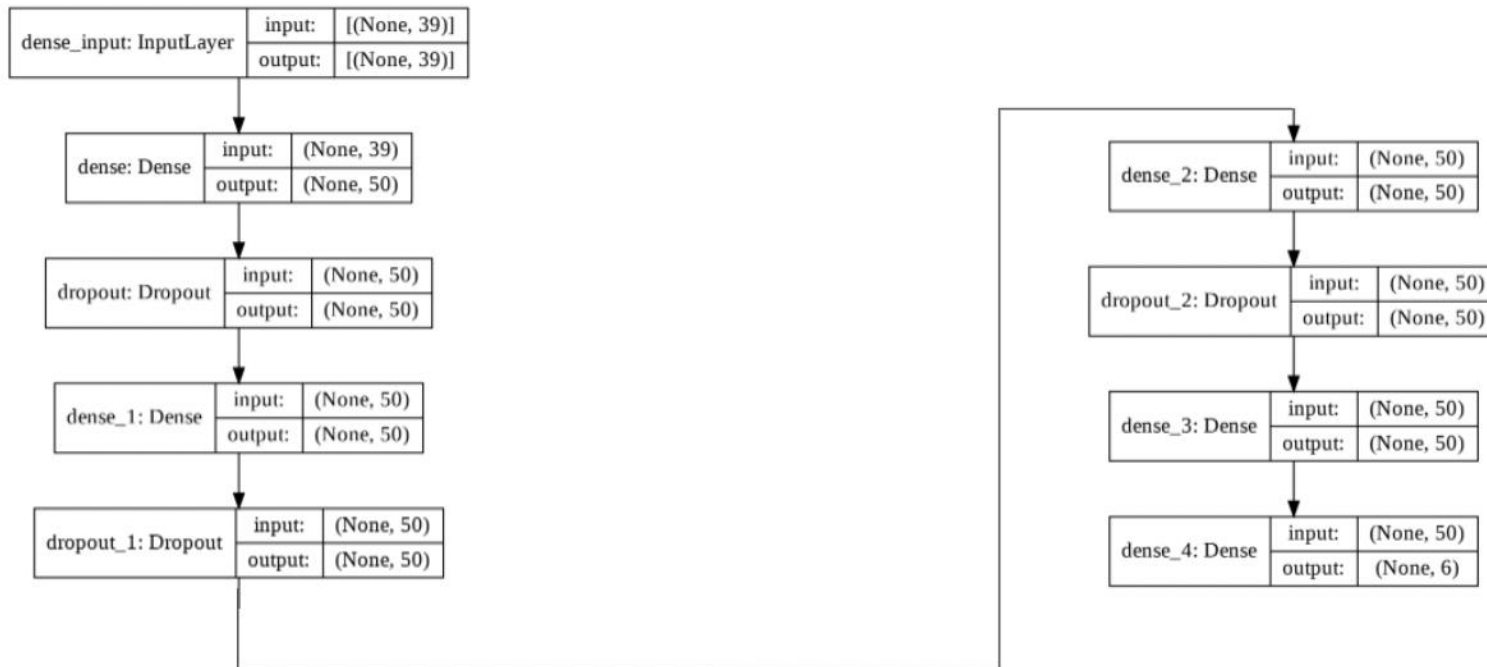
# Clustering accuracy for different number of neighbors on training meta data



# Neural network on the meta data

## Model architecture

- Tensorflow Keras sequential model

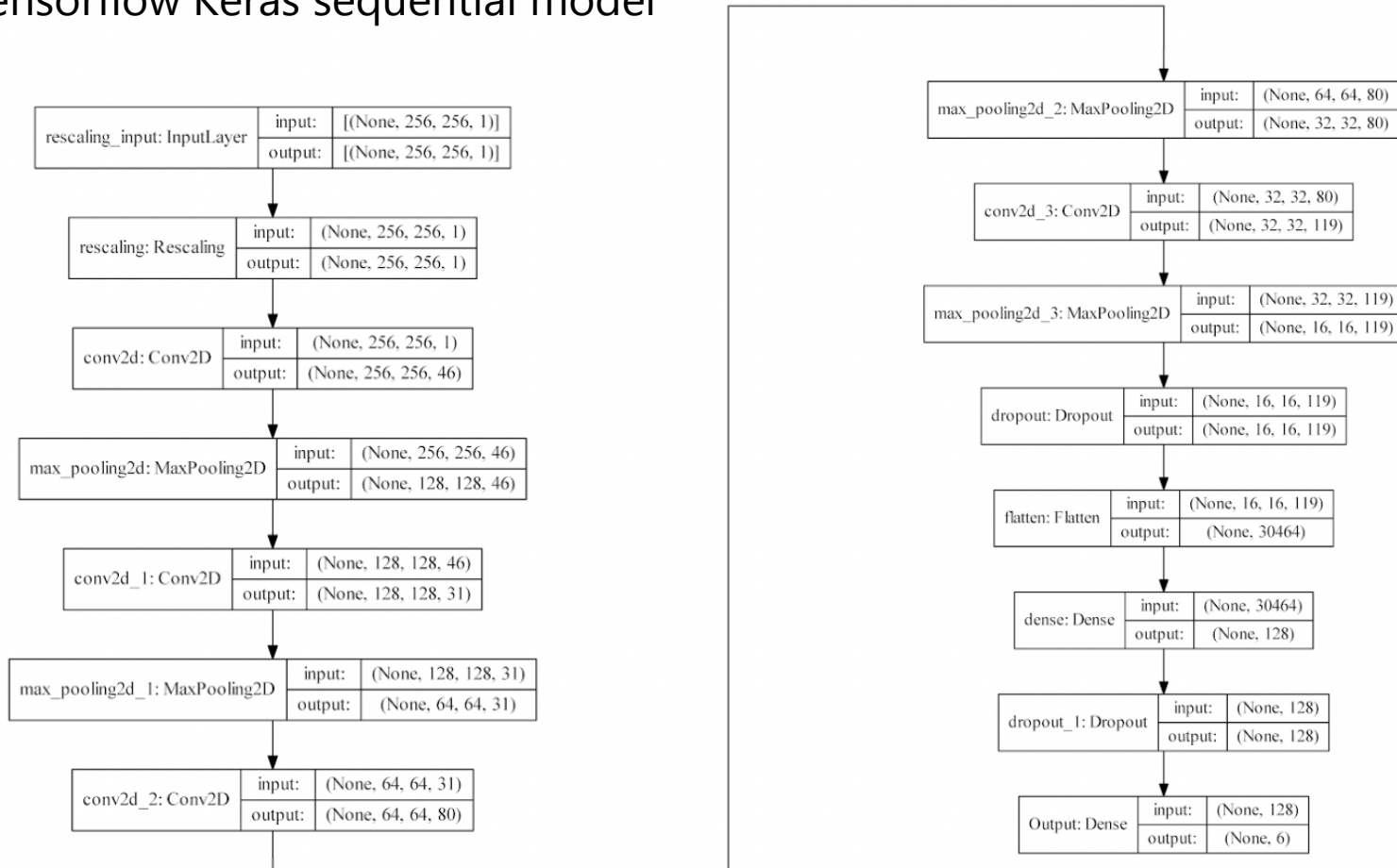




# Convolutional Neural Network on images

## Model architecture

- Tensorflow Keras sequential model

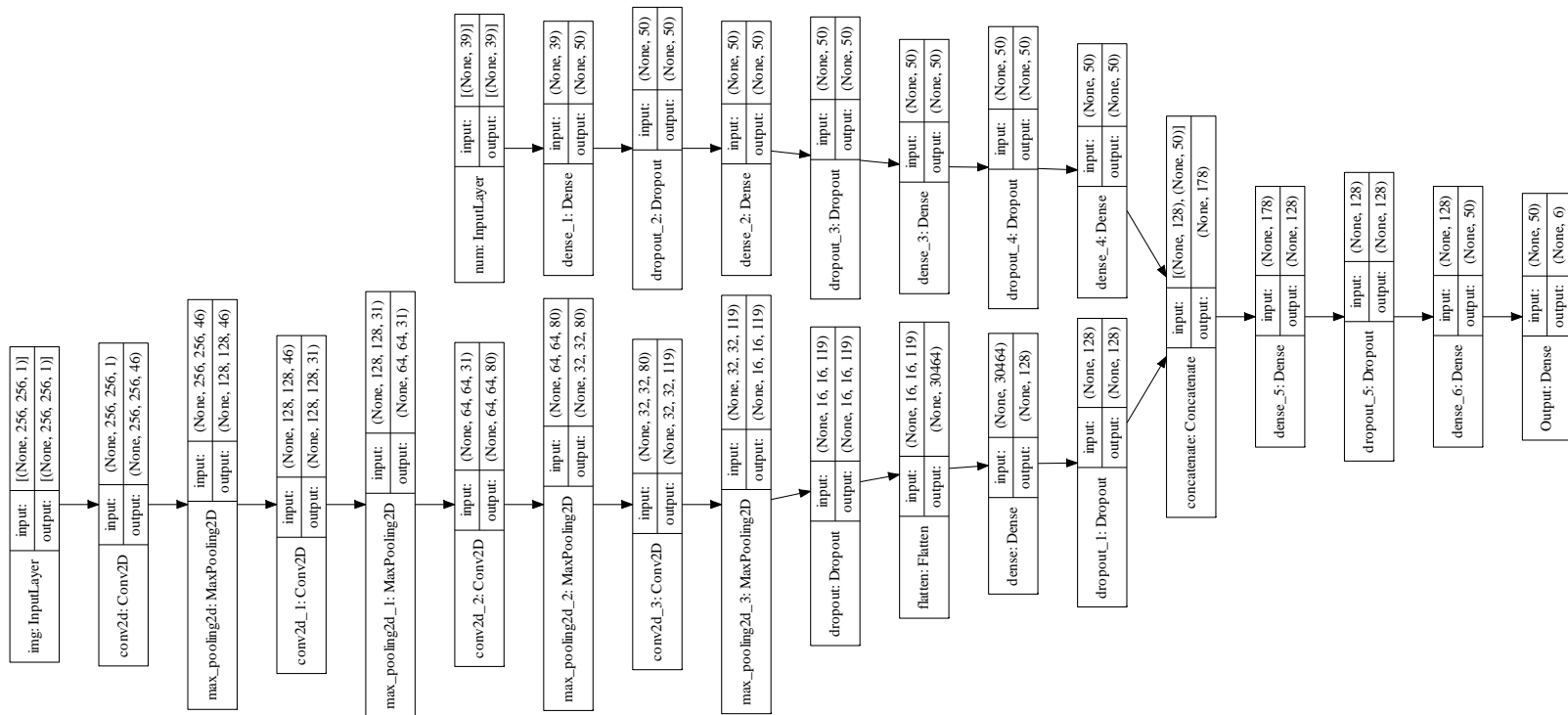


# Bayesian optimisation on CNN model

- Package used: Bayesian-optimization 1.2.0
- Parameters tuned
  - Number of filters for each Conv2D layer ( $n_1, n_2, n_3, n_4$ )
  - Drop out rate for each Dropout layer ( $\text{dropout1\_rate}, \text{dropout\_rate}$ )
  - l2 kernel regularization factor for Conv2D layers ( $\text{kernel\_reg\_rate}$ )
- Parameter space
  - $n_1$ : (16, 64),  $n_2$ : (16,64),  $n_3$ : (32, 128),  $n_4$ : (32, 128)
  - $\text{dropout1\_rate}$ : (0.1, 0.5),  $\text{dropout2\_rate}$ : (0.1, 0.5)
  - $\text{kernel\_reg\_rate}$ : (0.0001, 0.001)

# Mixed model Model architecture

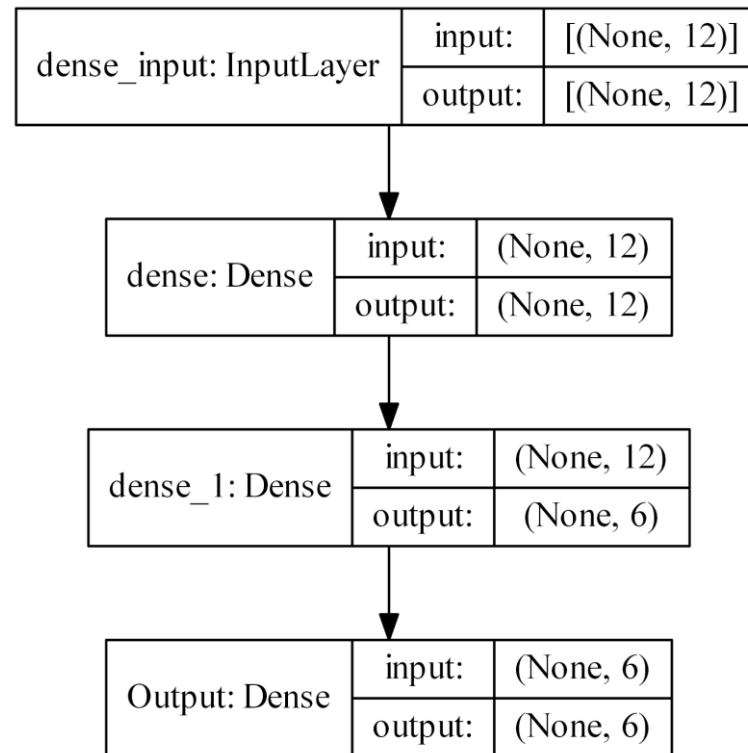
- Tensorflow Keras functional API model



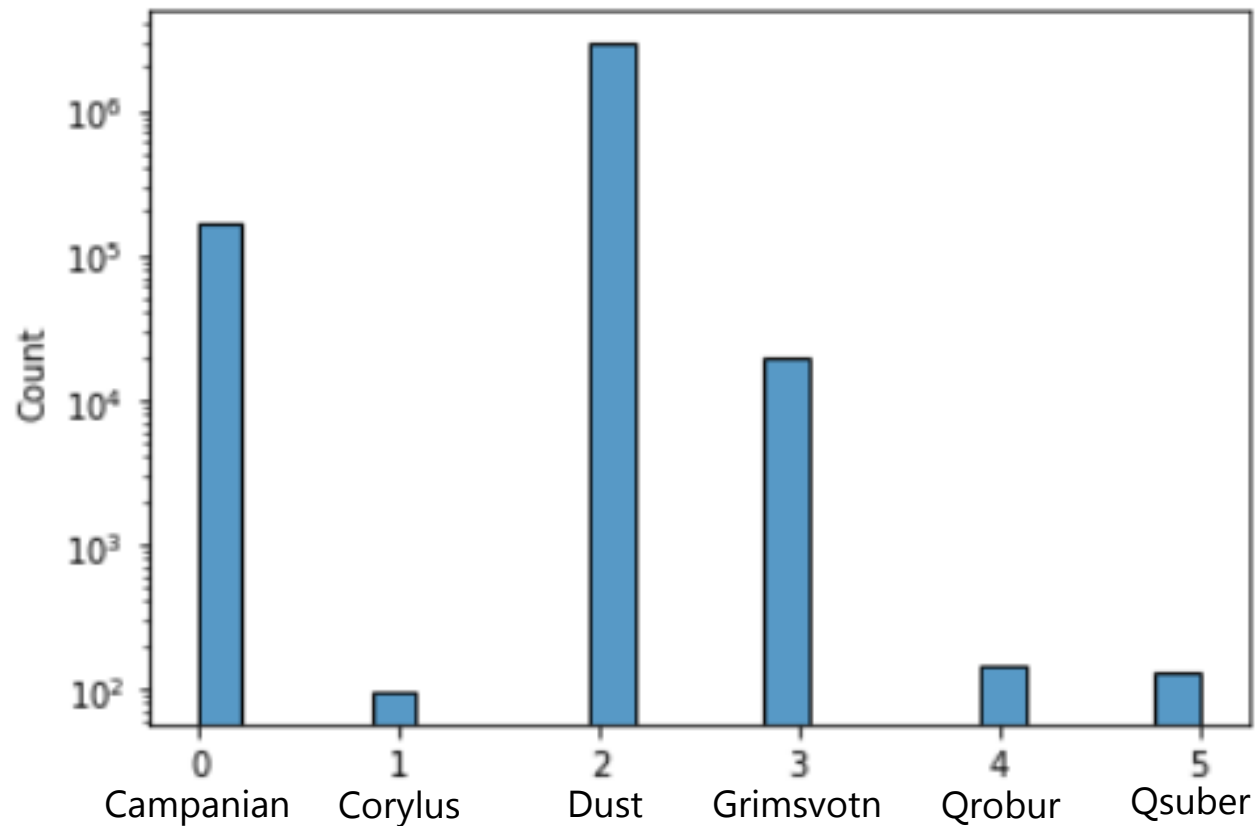
# Ensemble model

## Model architecture

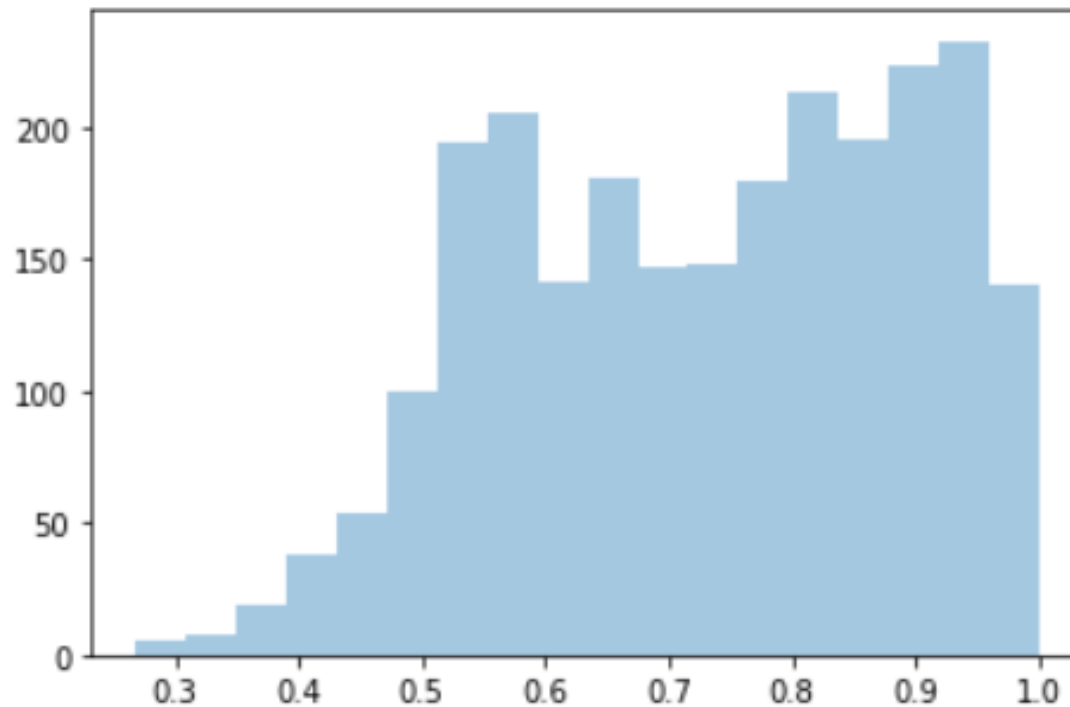
- Tensorflow Keras sequential model



# Log histplot of predictions on test ice core meta data

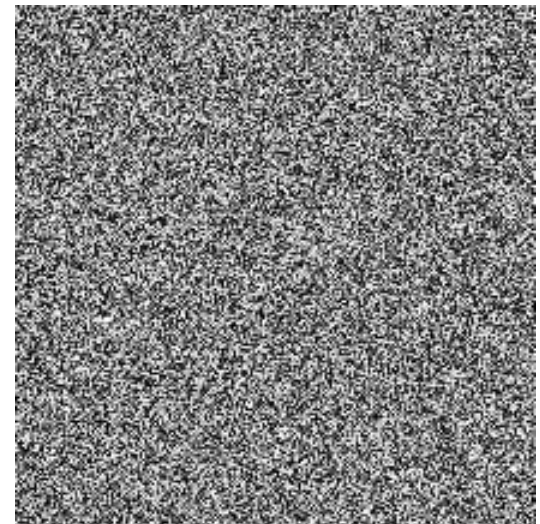
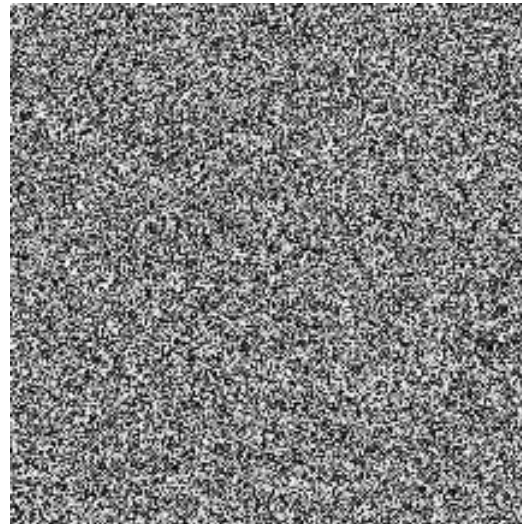
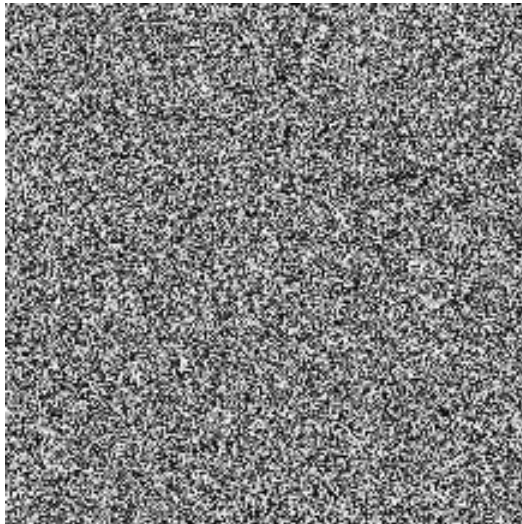


# Confidence level on wrong predictions from meta data



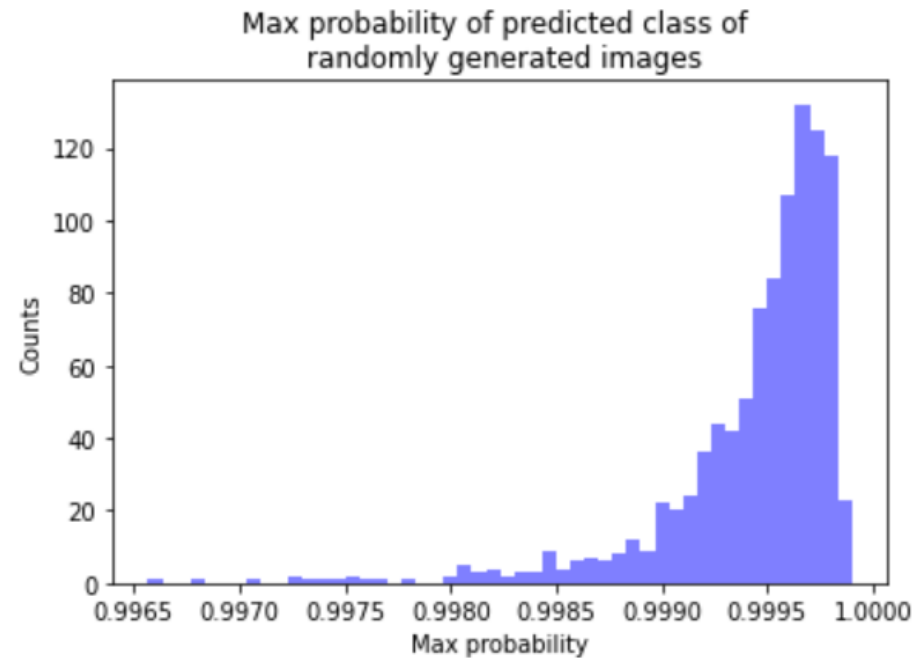
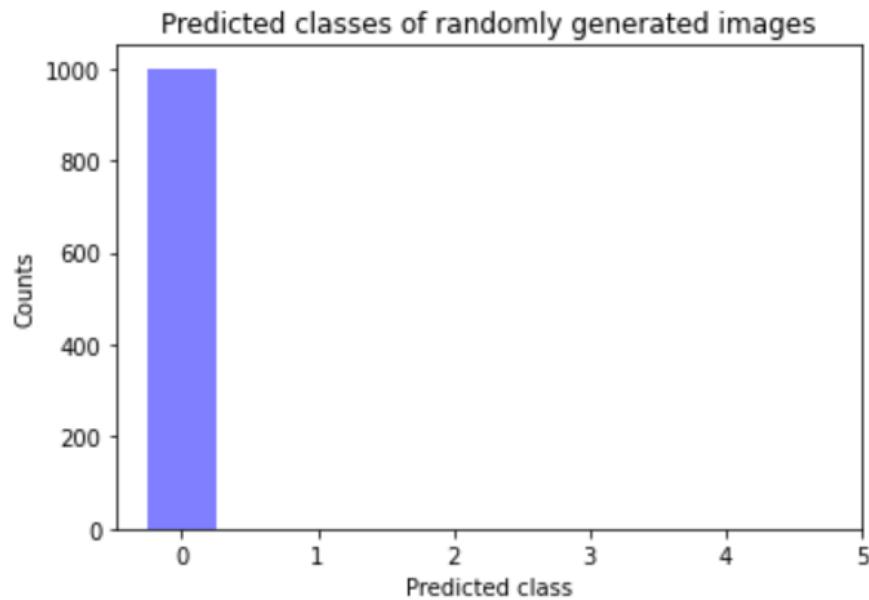
# Predictions on randomly generated pictures

- 1000 randomly generated images run through CNN to predict classes.



# Predictions on randomly generated pictures

- All randomly generated images are classified as class 0 i.e. Campanian with high confidence by the CNN.





# Predictions on random real images

- 8 random real images run through CNN to predict classes.



Predicted class: Grimsvotn  
Prediction confidence: 100%



Predicted class: Grimsvotn  
Prediction confidence: 100%



Predicted class: Grimsvotn  
Prediction confidence: 100%



Predicted class: Grimsvotn  
Prediction confidence: 100%



Predicted class: Campanian  
Prediction confidence: 79%



Predicted class: Campanian  
Prediction confidence: 99%

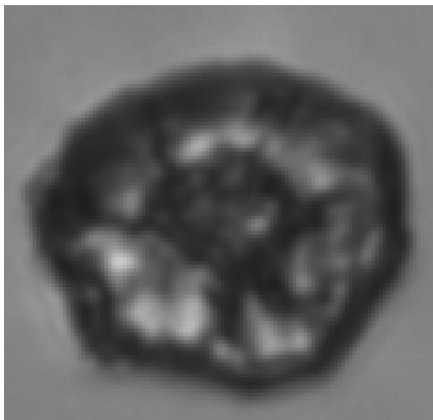


Predicted class: Grimsvotn  
Prediction confidence: 100%

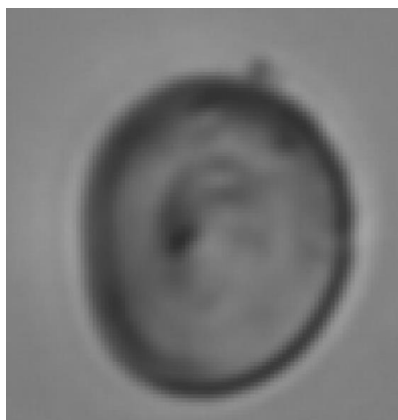


Predicted class: Grimsvotn  
Prediction confidence: 84%

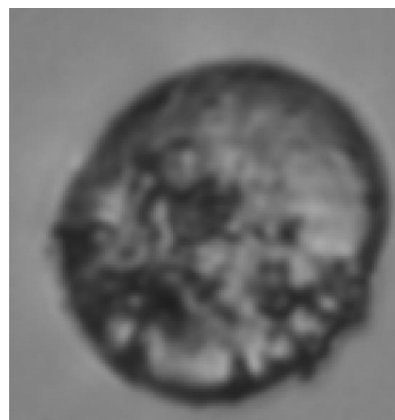
# Predictions on 6 most interesting images



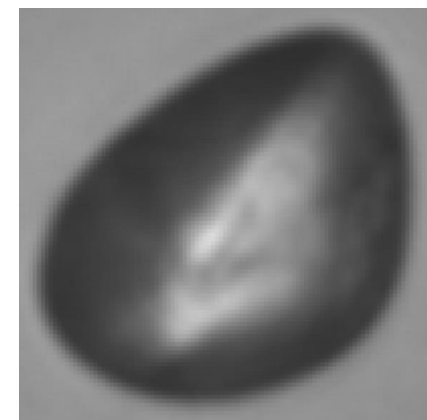
Predicted class: Corylus  
Prediction confidence: 100%



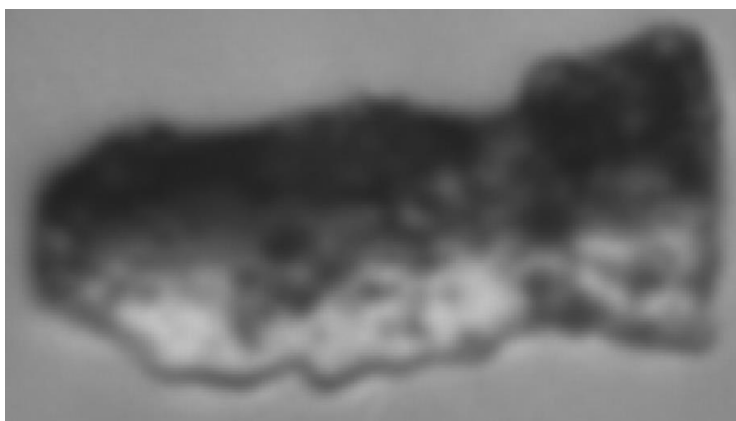
Predicted class: Qrobur  
Prediction confidence: 99%



Predicted class: Qsuber  
Prediction confidence: 62%



Predicted class: Corylus  
Prediction confidence: 100%



Predicted class: Corylus  
Prediction confidence: 45%



Predicted class: Qrobur  
Prediction confidence: 81%

# How much pollen do we expect the model to predict?

- Back of envelope-calculation

observed\_confusion \* ash\_frac\_in\_real\_data = expected\_frac\_pollen\_prediction  $\Rightarrow$

$(0.09+0.07+0.05)\% * 5\% + (0.01+0.03+0.06)\% * 5\% \sim 0.016\%$

