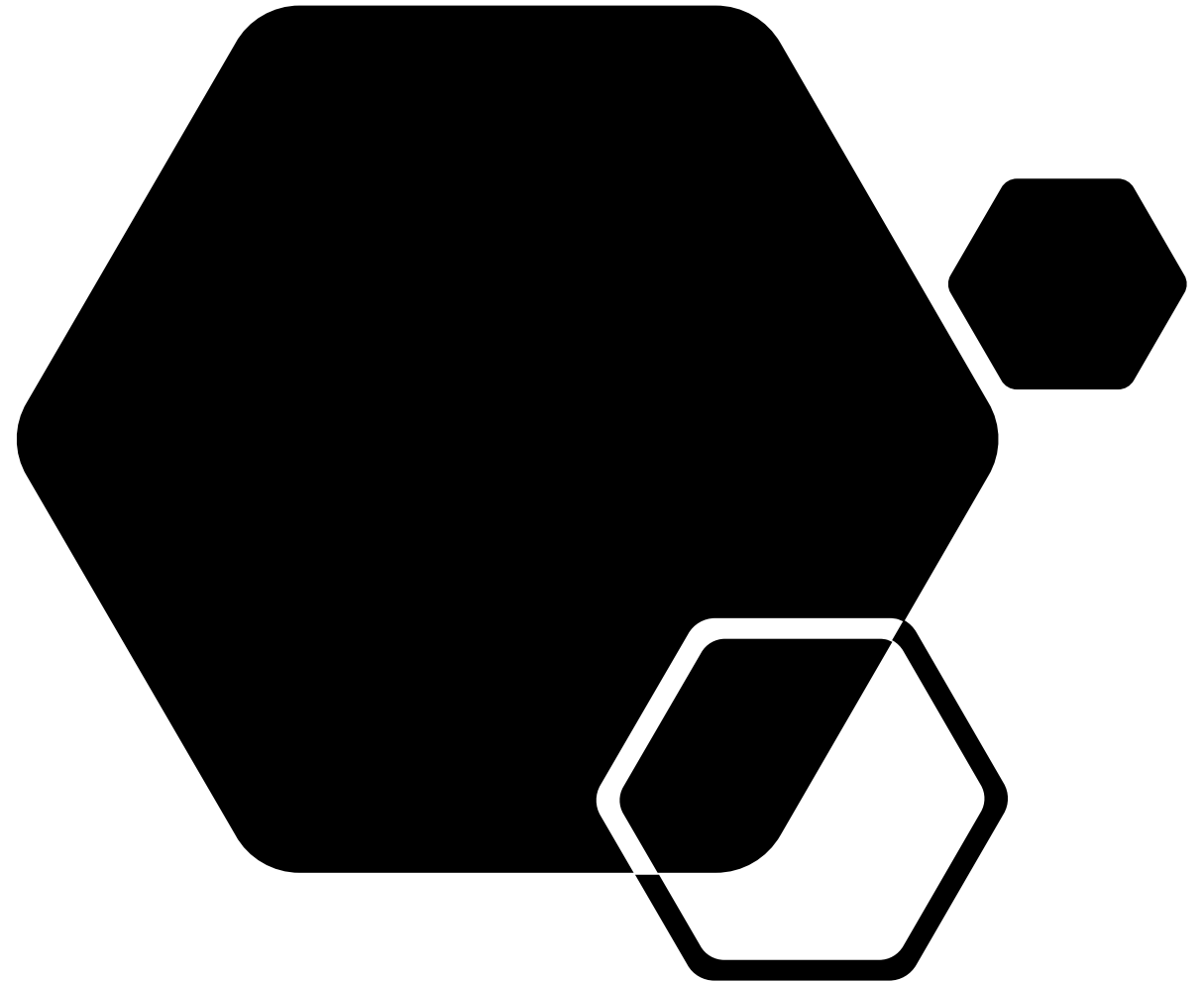


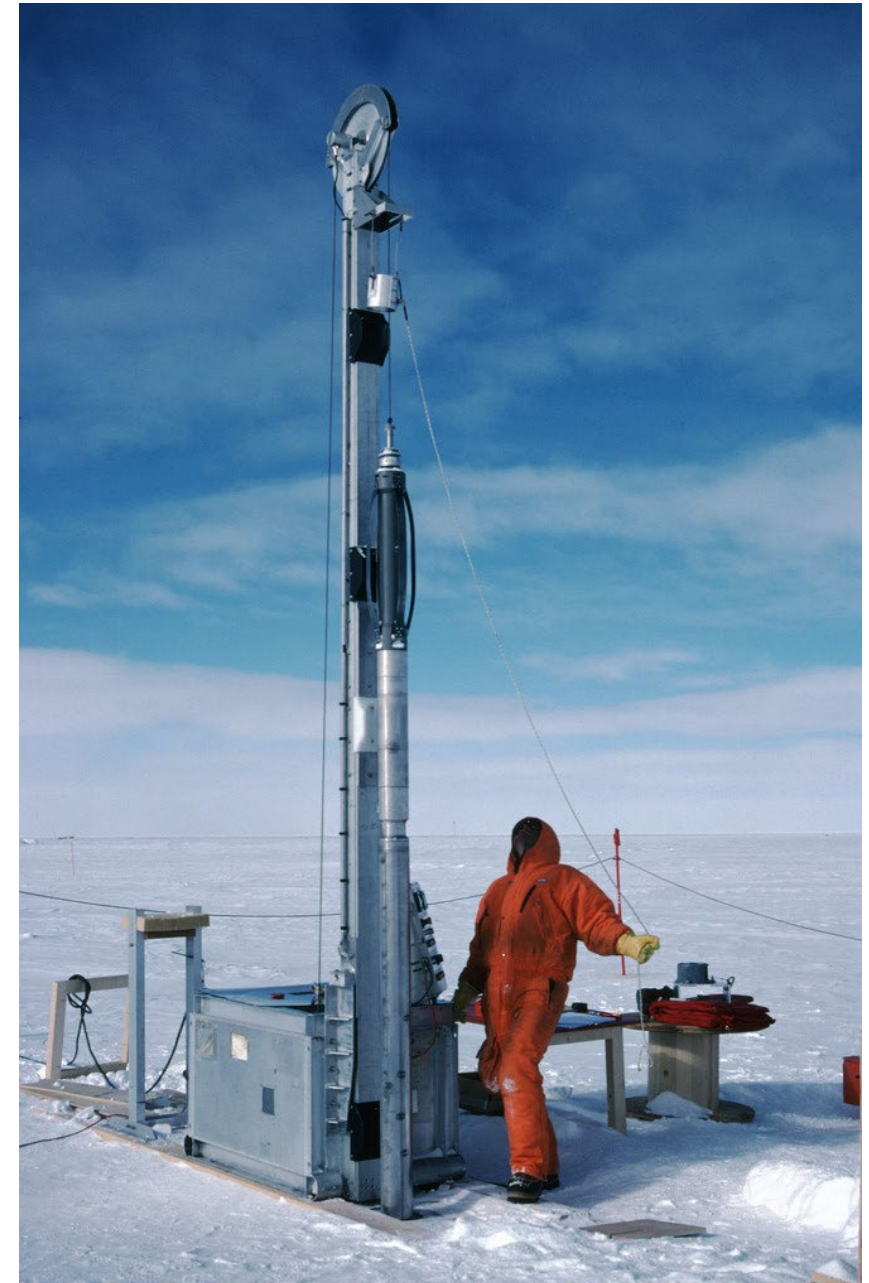
Ice Core Data Analysis

```
> get_participants()  
{  
  'Fabian Depenau Bjørnholt Jacobsen',  
  'Povl Filip Sonne-Frederiksen',  
  'Magnus Guldbæk Hansen',  
  'Marcus Melhedegaard Thomsen'  
}
```



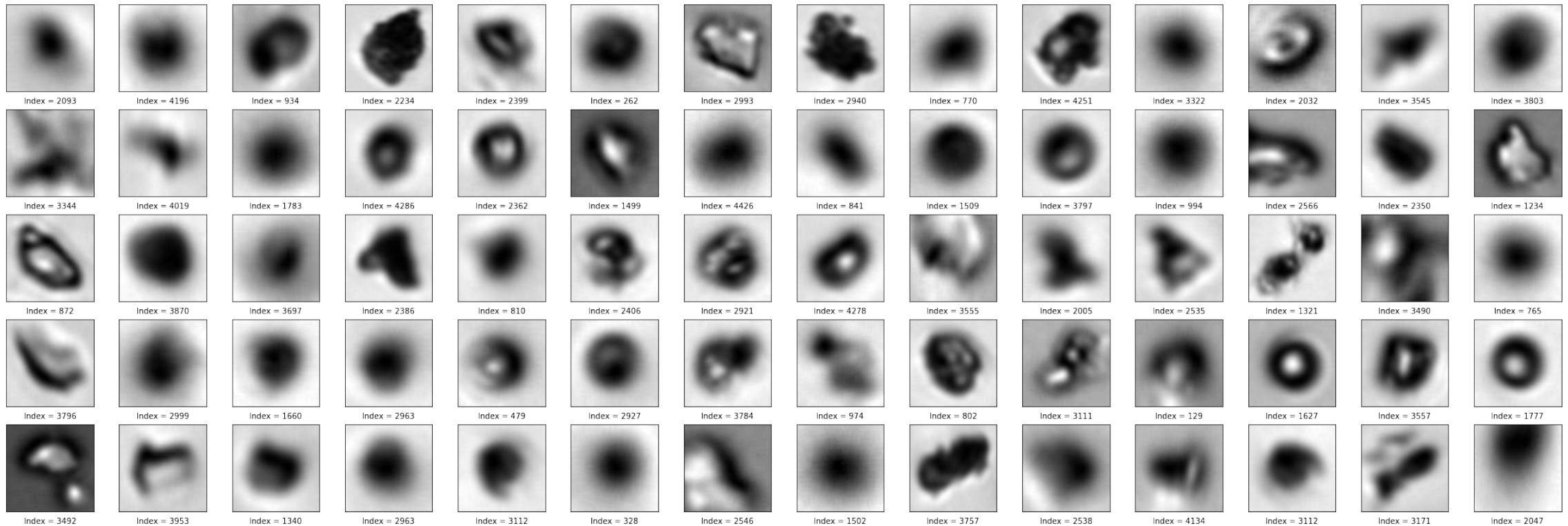
Overview

- Introduction
- Data
 - MNIST
 - Artificial dataset (labeled dataset from Nicolo)
 - Peruvian Ice Core Samples
- Methods
 - Classifier (Resnet18)
 - AutoEncoder (Resnet18 + Decoder)
 - Variational AutoEncoder (Resnet18 + Decoder)
- Analysis
 - Classifier
 - AutoEncoder
 - Variational AutoEncoder
- Findings
- Further Work



Introduction

The goal is to use dimensionality reduction to identify "interesting" objects within the Peruvian ice core dataset.



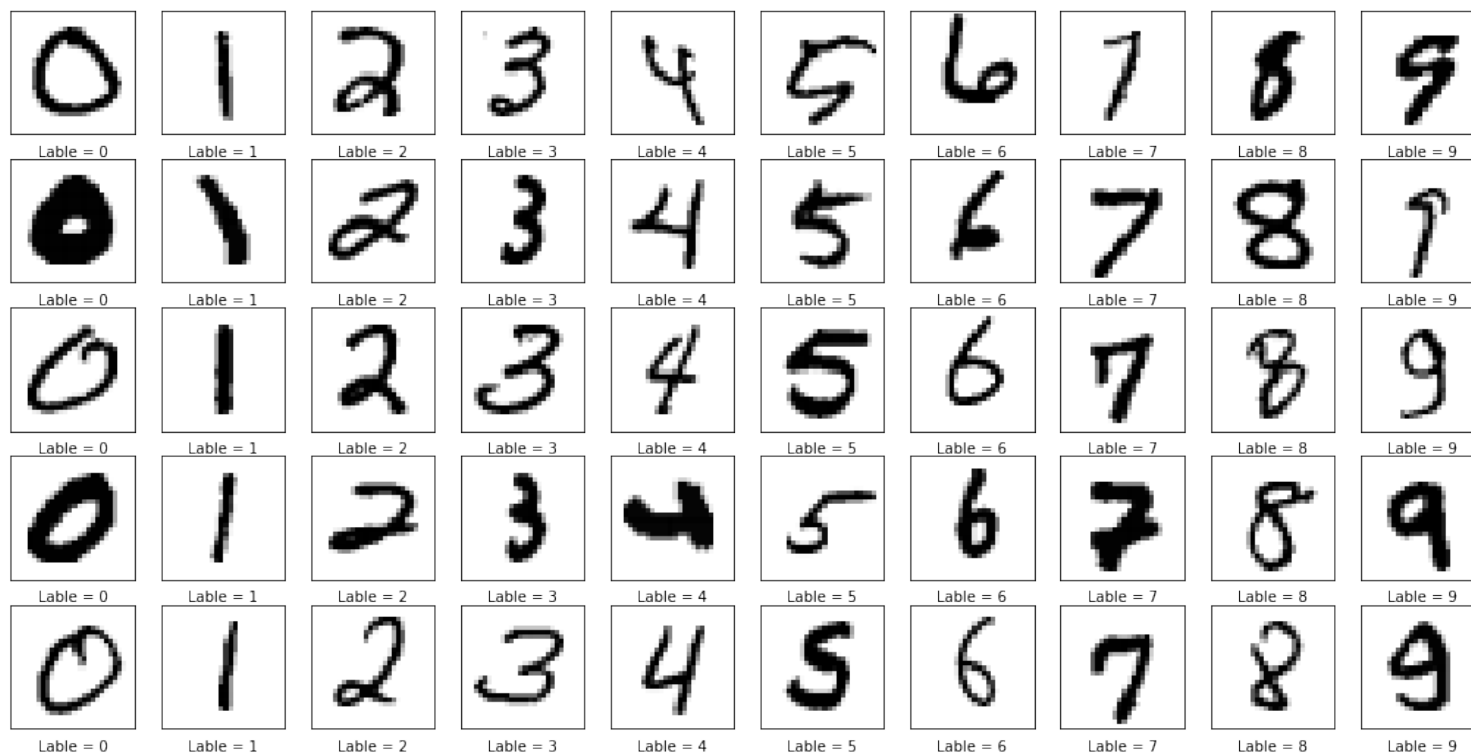
Data

A look at the datasets and their structures

MNIST

Handwritten digits

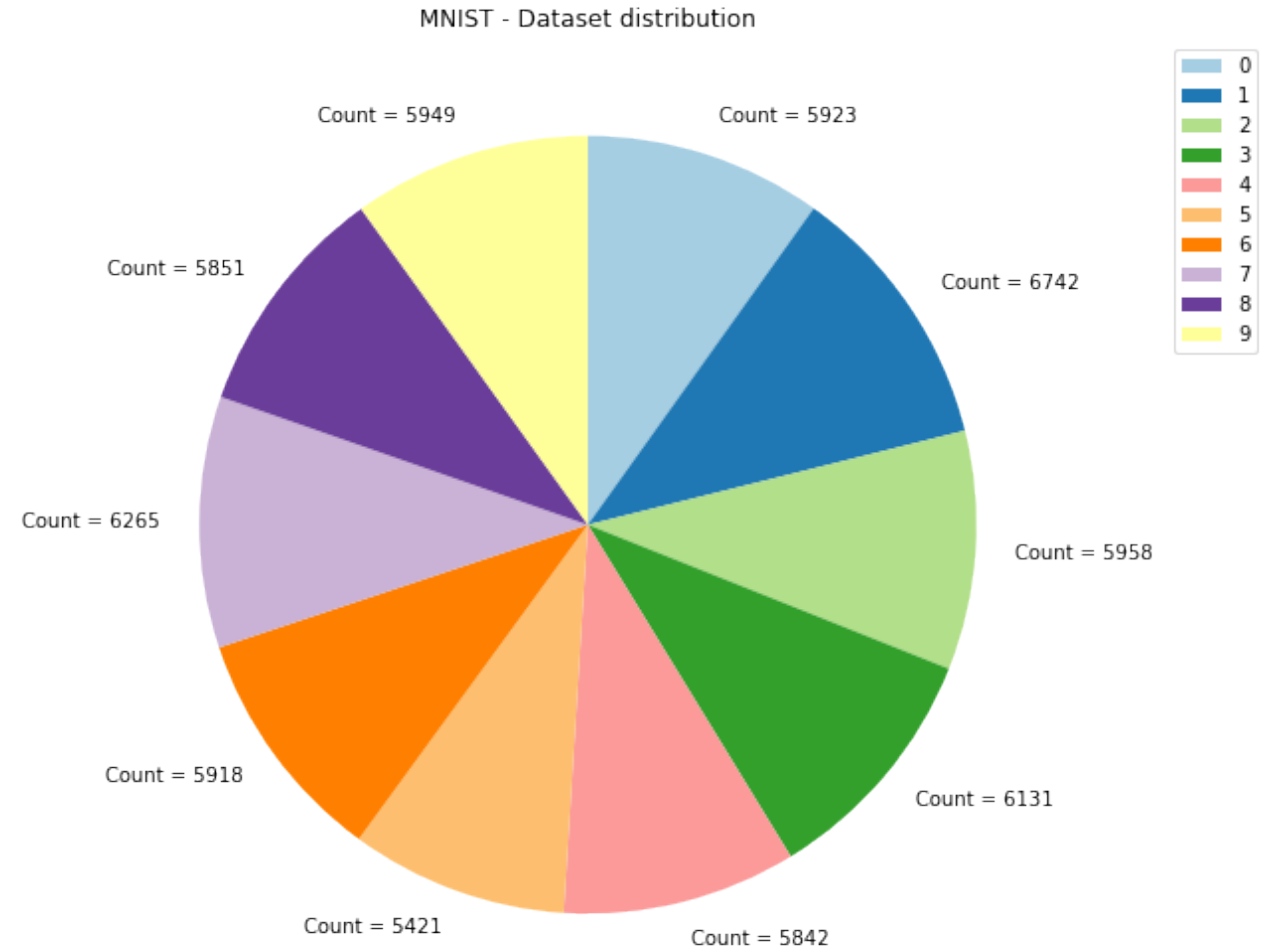
- 70.000 samples
- 28x28p size
- Balanced
- Preprocessed



MNIST

Handwritten digits

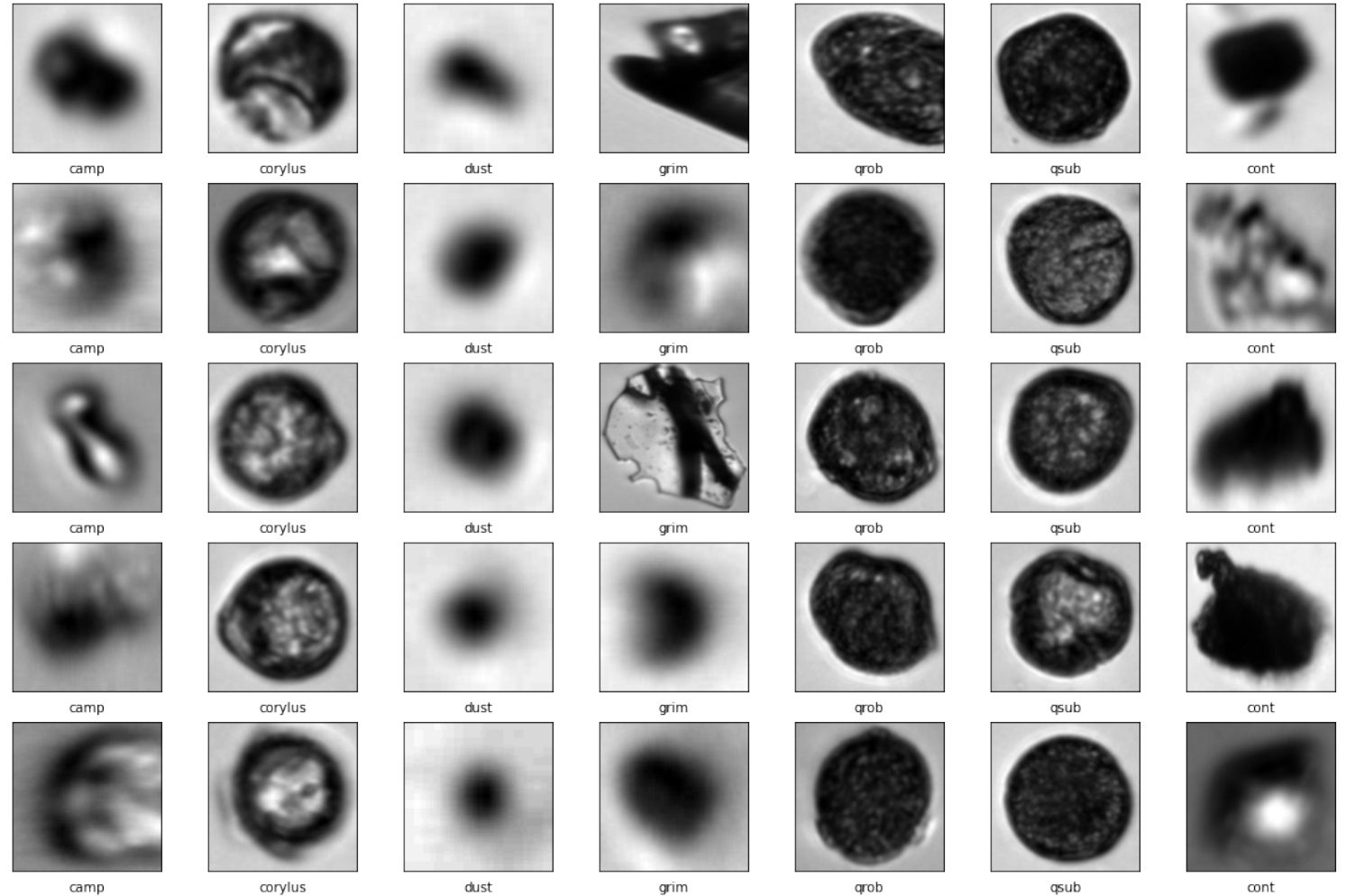
- 70.000 samples
- 28x28p size
- Balanced
- Preprocessed



Artificial

Image samples

- 145.242 samples
- Resized to 128x128p
- 34 scalars
- 7 Labels
- Unbalanced



Artificial

Image samples

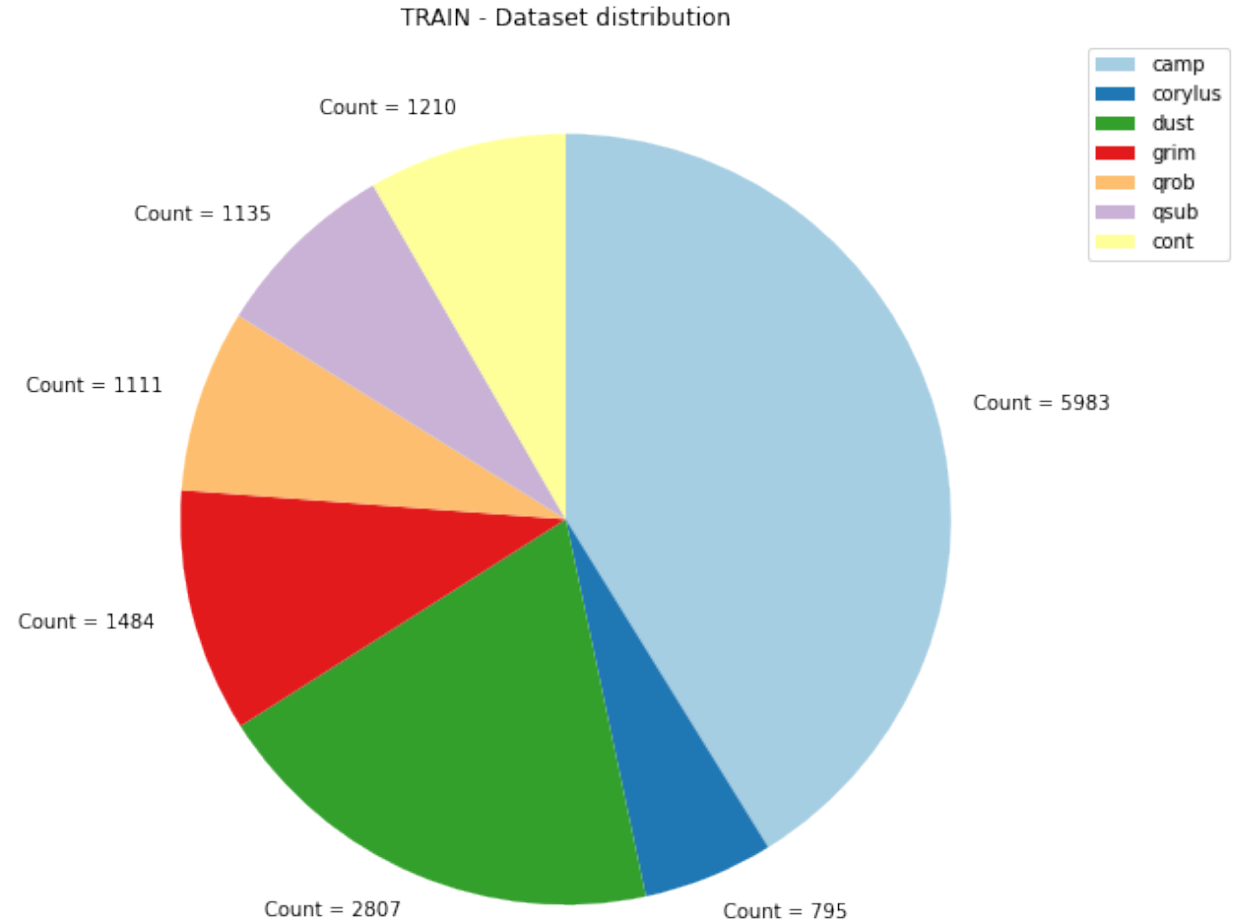
- 145.242 samples
- Resized to 128x128p
- 34 scalars
- 7 Labels
- Unbalanced

	Particle ID	Area (ABD)	Area (Filled)	Aspect Ratio	Biovolume (Cylinder)	...	Fiber Straightness	Fiber Score	Geodesic Aspect Ratio	Geodesic Length	Geodesic Thickness
0	3733	79.16	79.97	0.88	446.21	..	0.69	0	0.32	17.89	5.64
1	3526	9.85	9.85	0.67	47.64	..	1.2	0	1.0	3.93	3.93
2	27	711.31	716.33	0.71	11121.26	..	1.05	0	0.6	36.07	21.52
3	2780	54.6	54.6	0.5	293.64	..	0.88	0	0.37	13.87	5.19
4	3510	41.53	43.93	0.77	209.93	..	0.69	0	0.34	13.33	4.48
5	4940	49.15	49.15	0.6	227.46	..	0.79	0	0.3	14.92	4.41
6	251	492.14	3453.44	0.5	30415.77	..	1.01	0	0.41	61.07	25.19
7	3730	20.15	20.15	0.8	102.0	..	0.85	0	0.56	7.43	4.2
8	359	34.46	34.44	0.42	154.61	..	0.89	0	0.29	13.16	3.87
9	3041	2.16	2.16	0.55	12.67	..	1.06	0	1.0	2.53	2.53
10	538	3010.61	3043.2	0.56	57890.21	..	0.67	0	0.17	115.95	23.26
11	985	983.0	983.0	0.24	11772.3	..	0.99	0	0.21	69.68	34.67
12	3630	6.57	6.57	0.59	32.51	..	1.26	0	1.0	3.46	3.46
13	3024	24.32	25.07	0.53	95.4	..	0.79	0	0.29	11.43	3.26
14	4421	34.97	36.12	0.54	178.42	..	0.92	0	0.4	11.23	4.5
15	24910	11286.34	12293.88	0.29	217302.8	..	0.5	0	0.04	368.22	22.07
16	9695	1156.39	1157.79	0.94	22830.78	..	0.63	0	0.28	71.32	20.19
17	61.1	125.32	125.32	0.29	693.85	..	1.06	0	0.23	25.77	5.86
18	114	466.39	308.6	0.89	9016.04	..	1.24	0	1.0	22.56	22.56
19	7737	1339.99	3474.94	0.46	36981.02	..	0.67	0	0.12	112.62	13.86
20	26	30.36	30.36	0.71	35.07	..	1.12	0	1.0	4.12	4.12
21	2304	27.5	27.5	0.44	20.27	..	0.6	0	0.12	10.07	2.22
22	82	494.23	201.66	0.87	9876.65	..	1.22	0	1.0	23.26	23.26
23	822	3541.71	3534.44	0.8	300645.08	..	0.71	0	0.26	124.11	32.16
24	3002	3.06	3.06	0.53	26.3	..	1.21	0	1.0	3.23	3.23
25	456	846.26	846.26	0.32	11550.66	..	1.05	0	0.29	36.47	36.14
26	1718	791.41	791.41	0.95	15907.62	..	1.23	0	1.0	27.26	27.26
27	4090	39.84	39.84	0.81	113.07	..	1.13	0	1.0	5.24	5.24
28	2251	21.55	21.55	0.69	124.98	..	1.24	0	1.0	5.42	5.42
29	2358	30.19	30.19	0.55	53.09	..	1.3	0	1.0	4.07	4.07
30	357	0.4	0.4	0.39	1.33	..	0.95	0	1.0	1.19	1.19
31	5305	815.77	815.77	0.76	13897.48	..	0.87	0	0.46	43.73	20.12
32	3592	24.01	24.01	0.24	27.0	..	0.77	0	0.03	37.65	0.96
33	386	3857.91	3907.26	0.63	34429.87	..	0.34	0	0.04	218.66	9.17
34	17161	401.0	402.07	0.51	4071.36	..	0.86	0	0.3	36.6	11.59
35	967	21.24	21.24	0.8	115.7	..	1.14	0	1.0	5.28	5.28
36	2859	3076.36	3083.77	0.86	21755.35	..	0.89	0	0.5	48.55	23.94
37	227	25.1	25.1	0.55	369.29	..	0.83	0	0.39	11.28	4.37
38	49	465.67	499.97	0.9	6801.59	..	0.82	0	0.48	33.64	36.04
39	827	1139.05	1143.45	0.81	15046.28	..	0.67	0	0.2	77.54	25.72
40	1382	919.66	942.76	0.82	15861.03	..	0.79	0	0.39	30.65	39.97
41	29398	80.1	80.1	0.6	388.58	..	0.72	0	0.24	20.6	4.9
42	4874	21.67	21.67	0.63	118.23	..	0.97	0	0.63	7.28	4.55
43	5582	958.78	958.78	0.89	20170.92	..	0.92	0	0.61	41.14	24.98
44	399	476.19	504.41	0.86	4743.64	..	0.56	0	0.21	50.98	30.88
45	3024	39.91	41.27	0.57	203.68	..	0.82	0	0.37	12.34	4.58
46	4216	47.71	48.29	0.2	135.1	..	0.8	0	0.11	24.38	2.66
47	3207	7.0	7.0	0.84	30.93	..	1.04	0	1.0	3.4	3.4
48	3247	602.16	604.12	0.6	4836.81	..	0.57	0	0.13	20.58	9.34
49	1371	394.89	394.89	0.55	3990.32	..	1.19	0	0.6	27.16	36.19

Artificial

Image samples

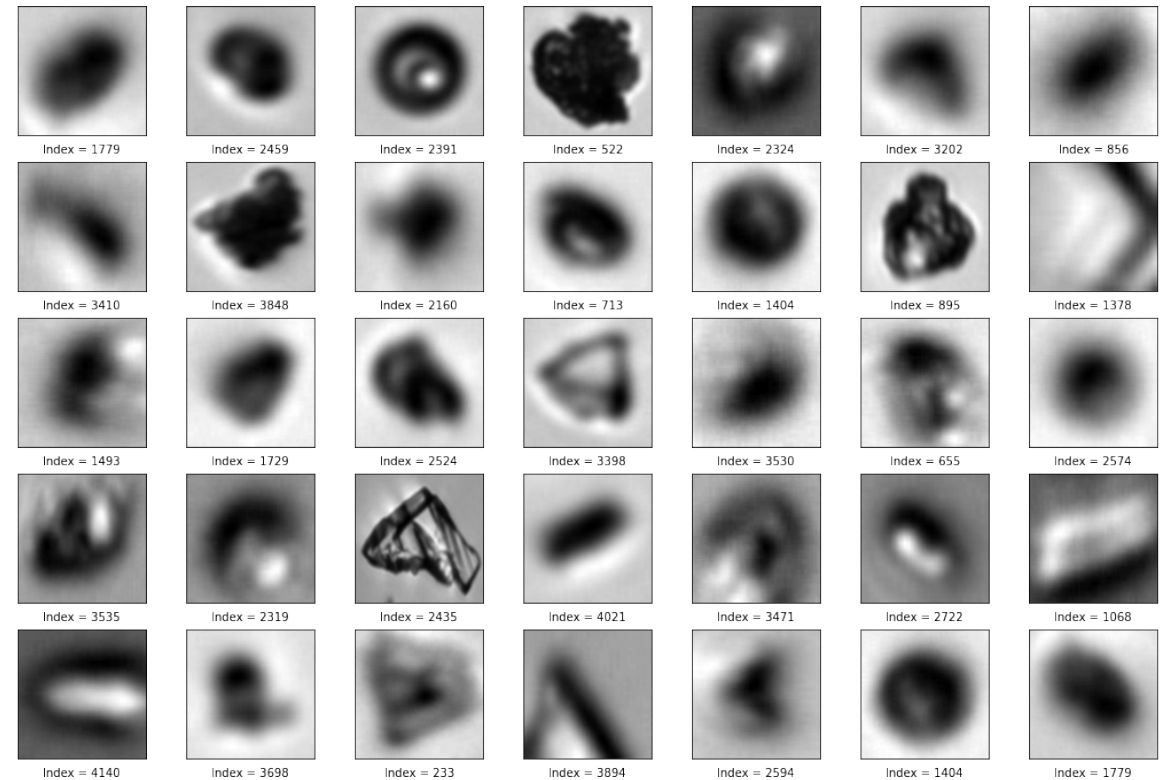
- 145.242 samples
- Resized to 128x128p size
- 34 scalars
- 7 Labels
- Unbalanced



Peruvian Ice Core Samples

Image samples

- 102.763 samples
- Resized to 128x128p
- 34 scalars
- Probably extremely unbalanced with a lot being dust

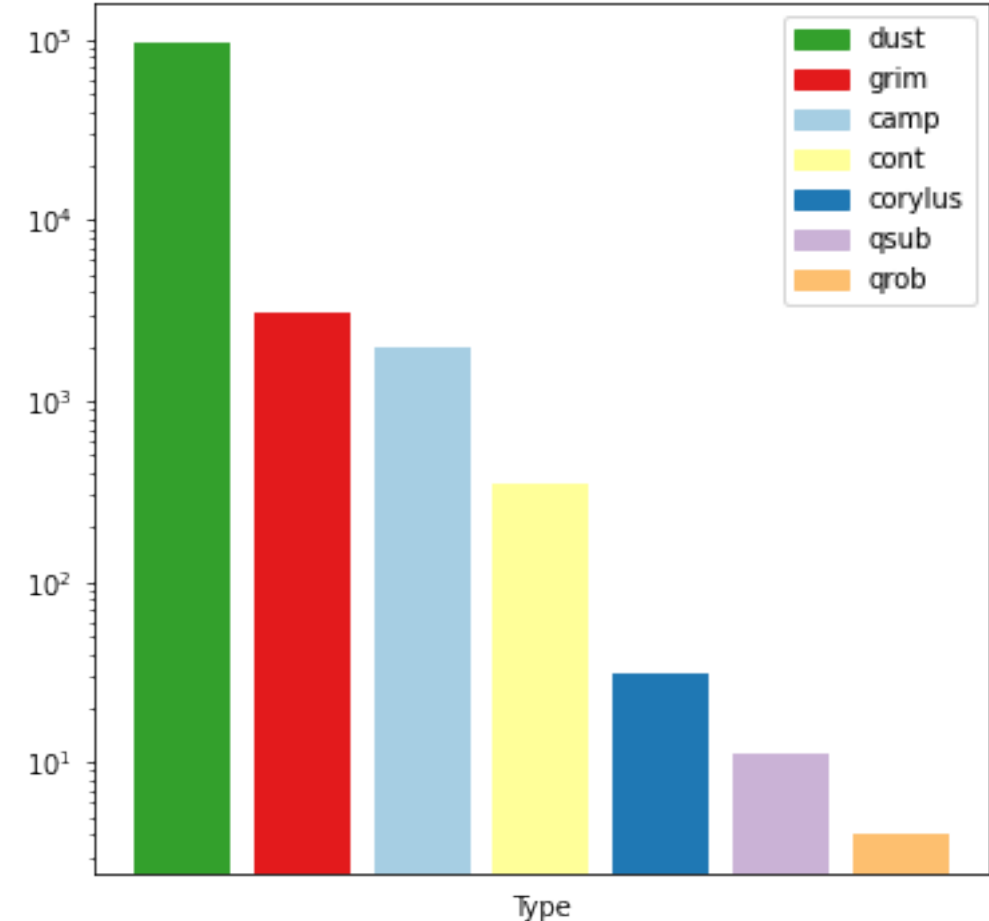


Peruvian Ice Core Samples

Image samples

- 102.763 samples
- Resized to 128x128p size
- 34 scalars
- Probably extremely unbalanced with a lot being dust

Presumed class distribution in the Peruvian Dataset



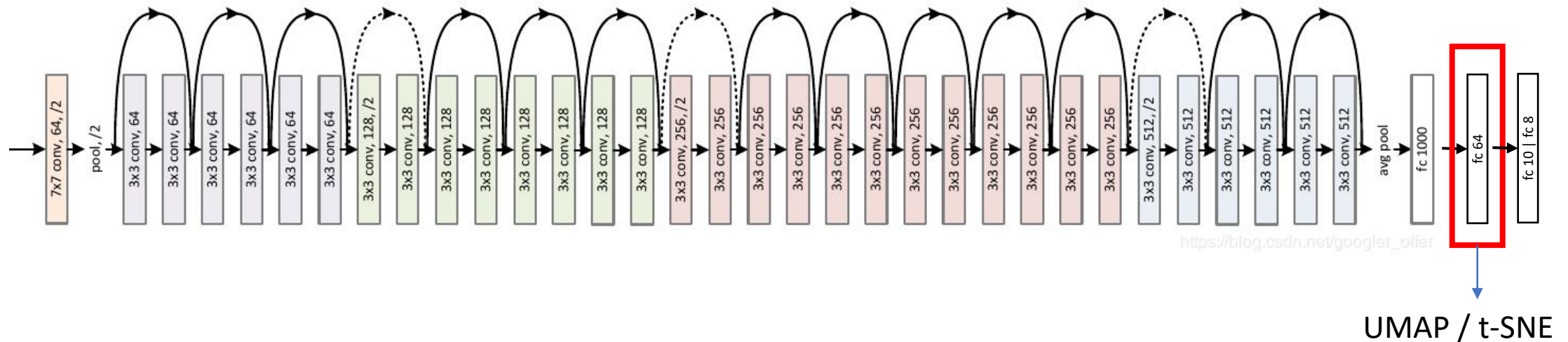
We remove all that the classifier characterizes as dust from the Peruvian set and look at the rest.
We include only images above 15x15 pixels in size

Methods

How did we engage with the data

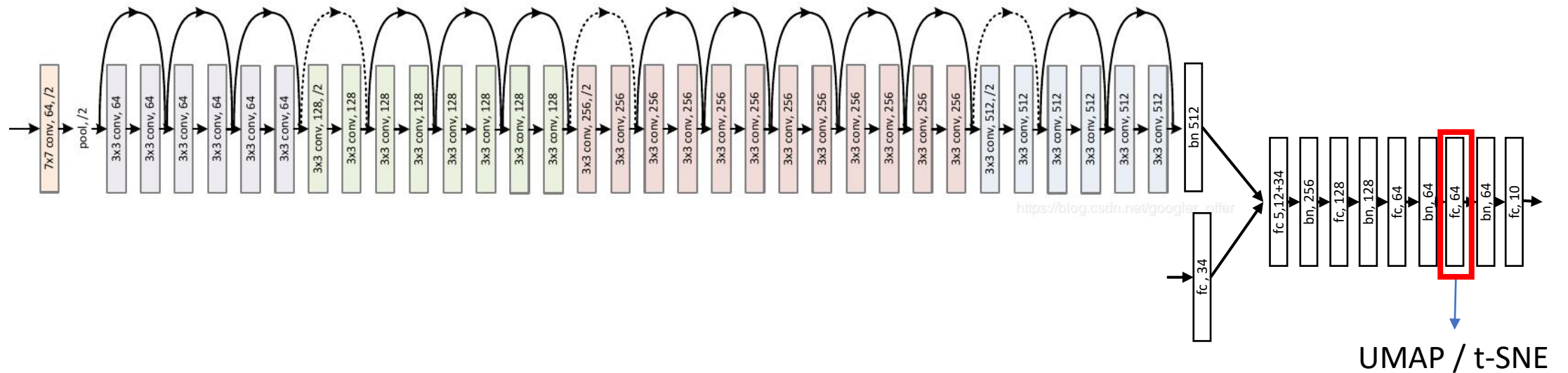
Classifier – Resnet18 on MNIST

- See if classification can be used to find unknow categories by mapping the n-1 layer in the network



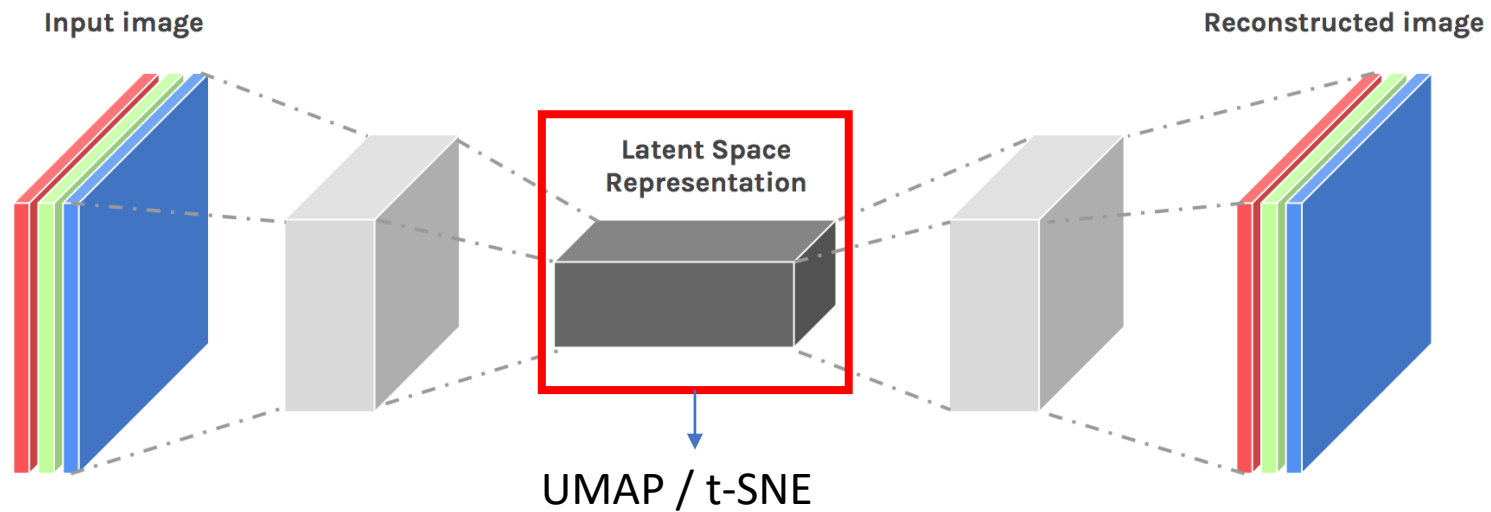
Classifier – Resnet18 on Artificial and Peruvian Datasets

- See if classification can be used to find unknown categories by mapping the n-1 layer in the network



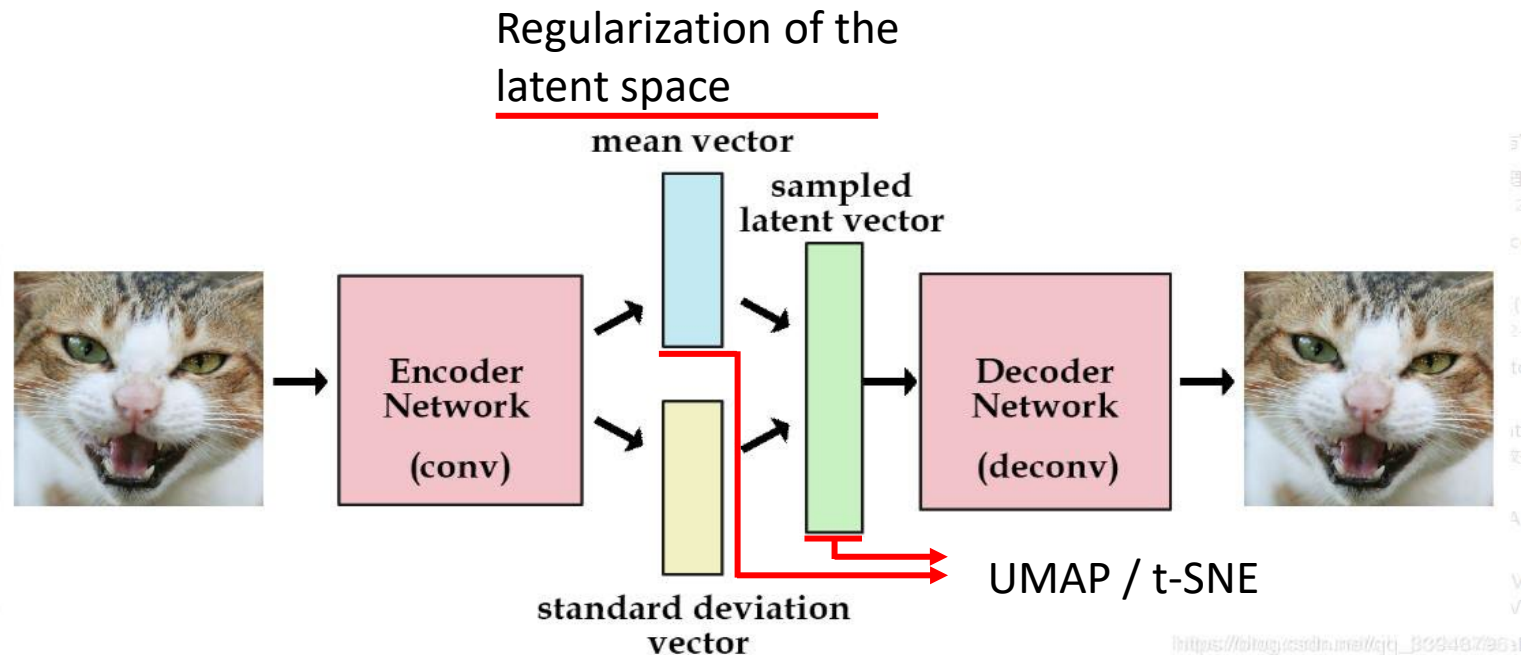
AutoEncoder

- See if the bias can be reduced by encoding a latent space

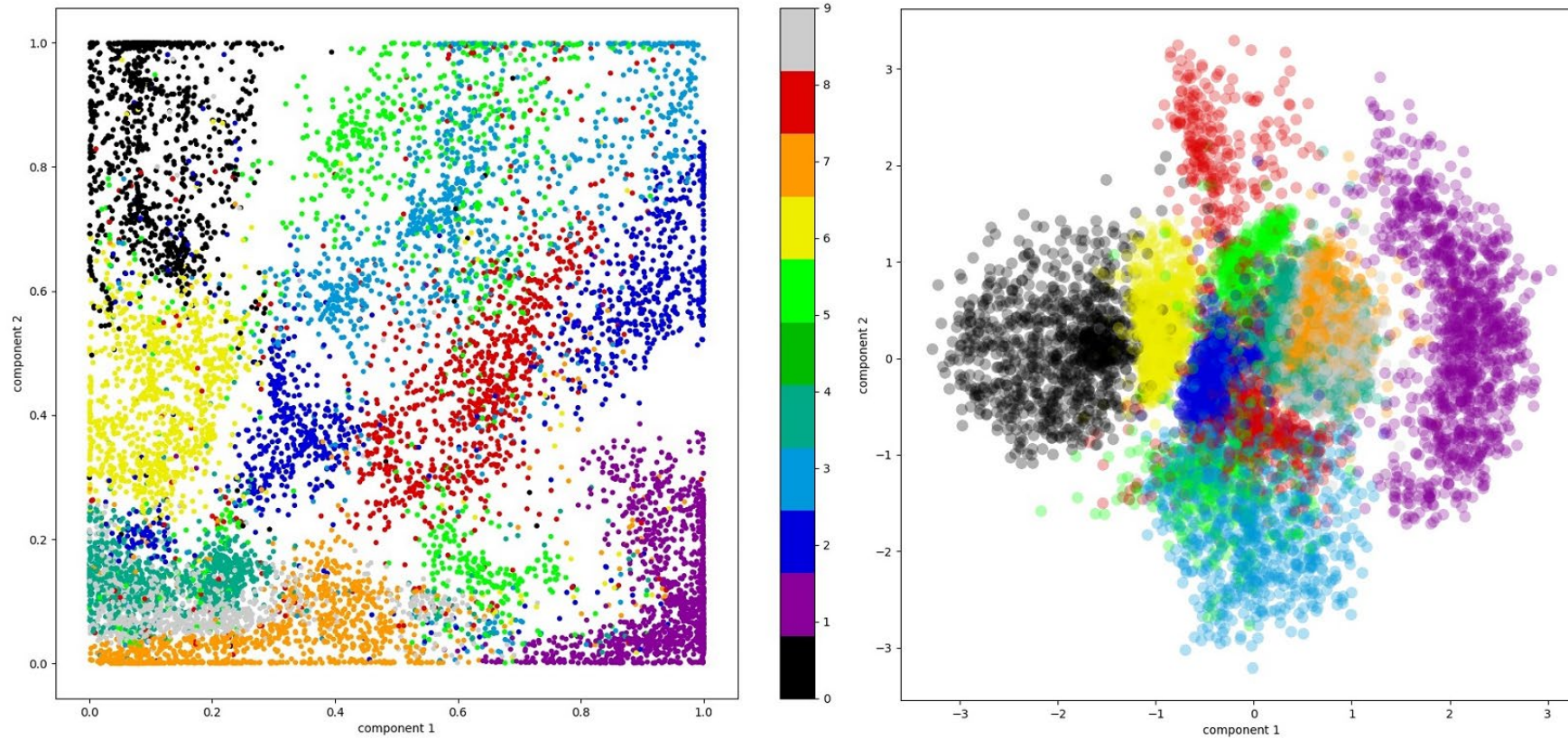


Variational AutoEncoder

- Regularizing the latent space to reduce overfitting



Example of latent space for AE vs VAE

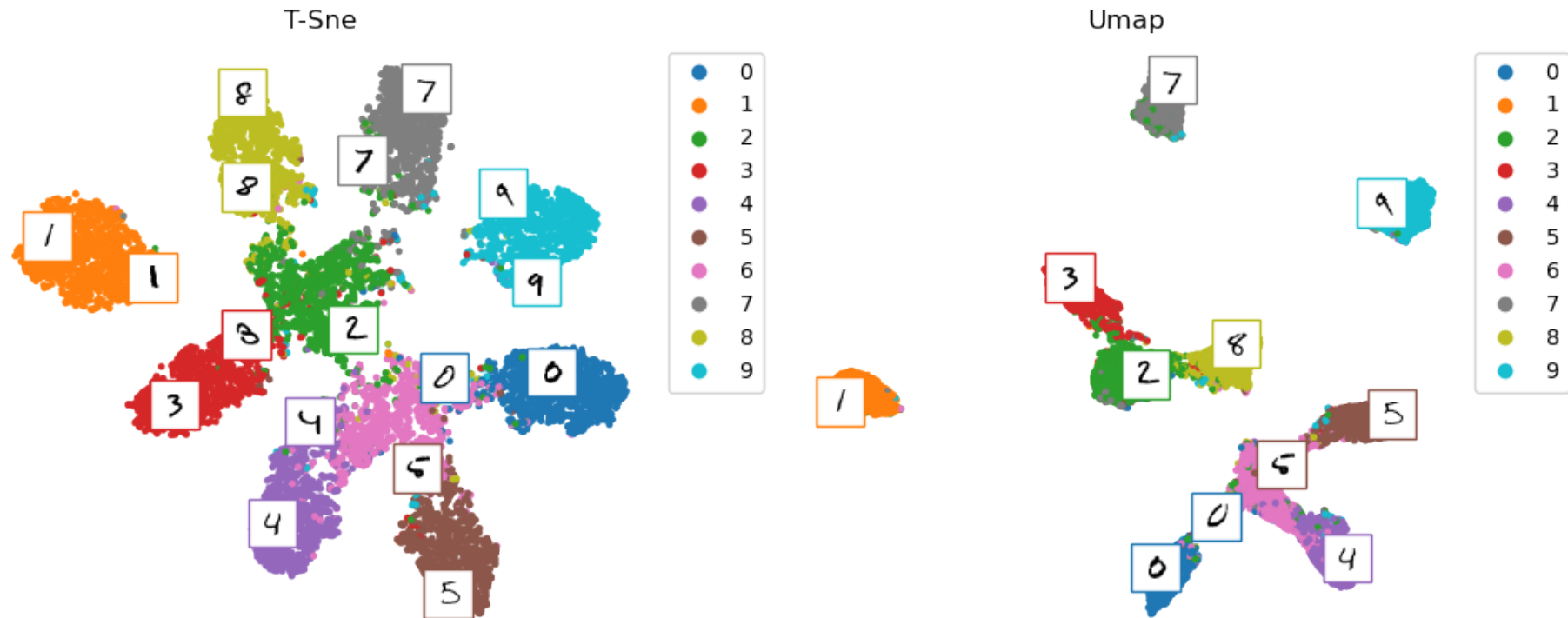


Analysis

Our results

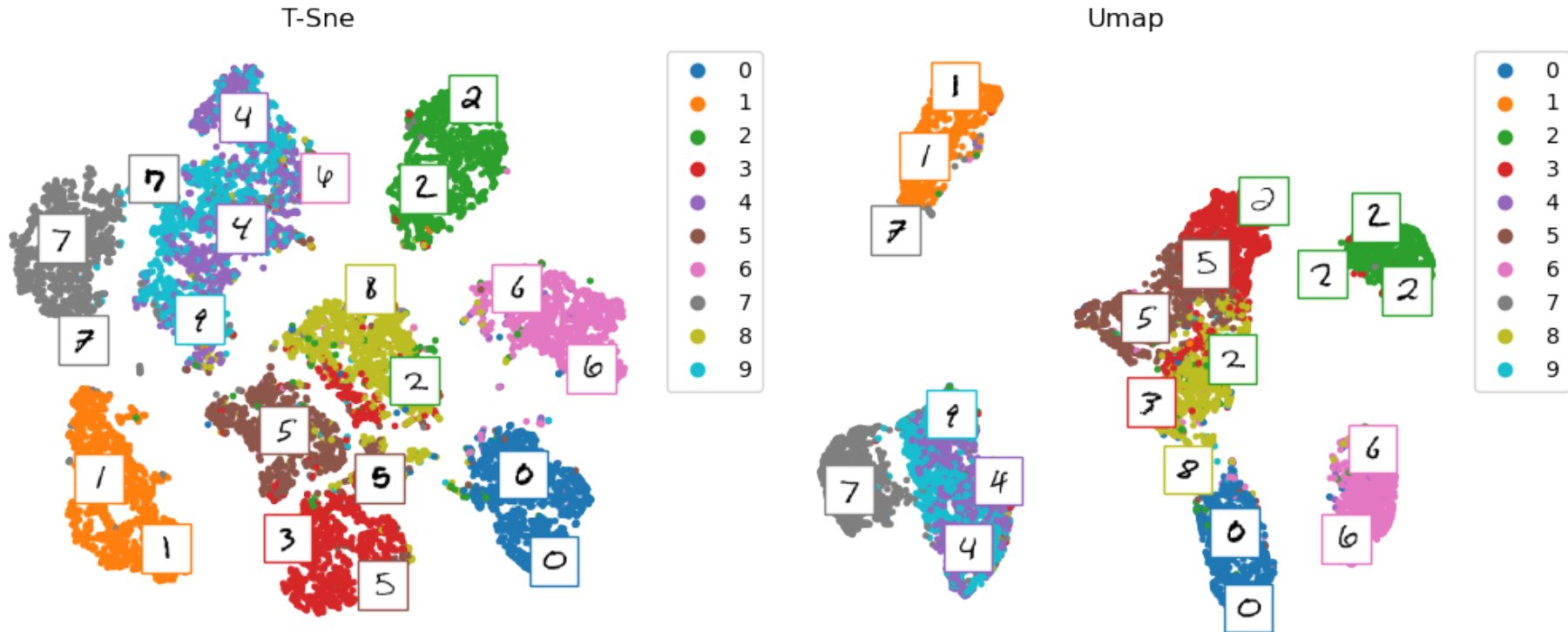
Classifier – MNIST

Mapping the second to last layer in ResNet18 where [2, 6] have been left out of the MNIST Train Dataset



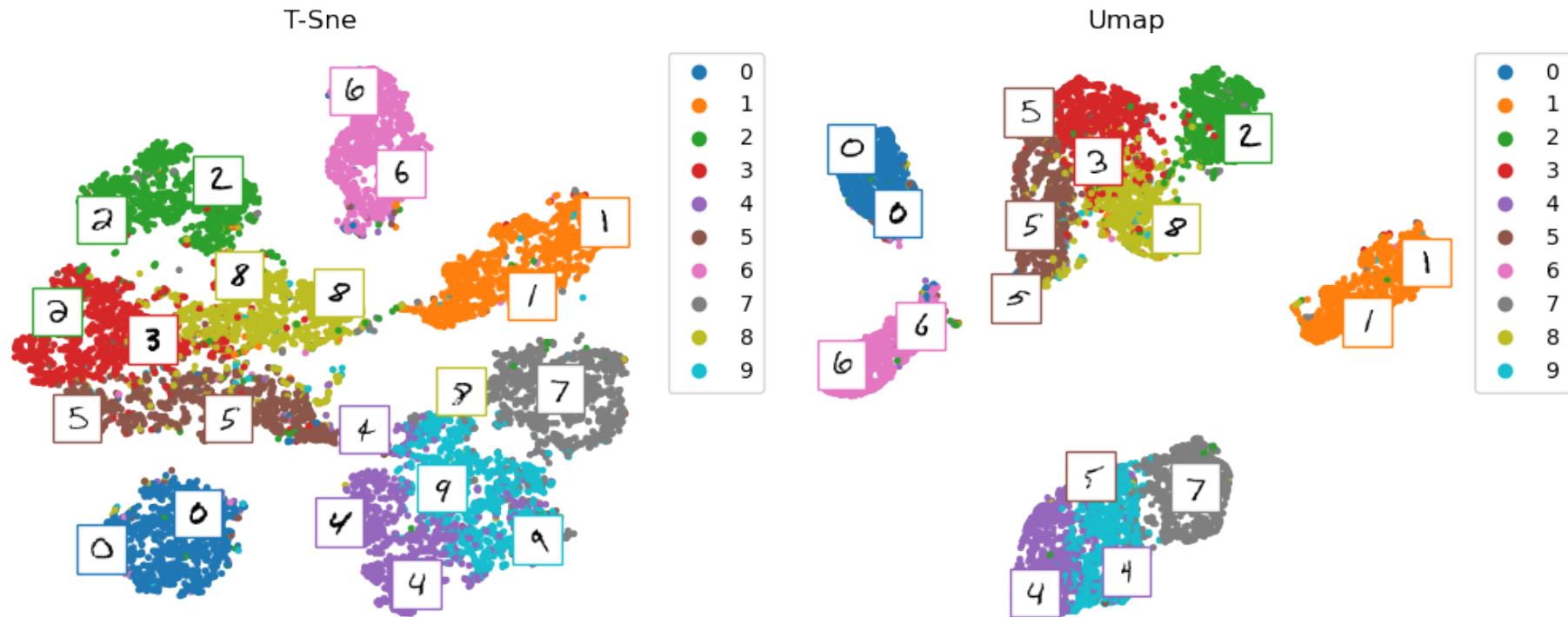
AutoEncoder – MNIST

Mapping the latent space of aM AutoEncoder trained on the MNIST Train Dataset

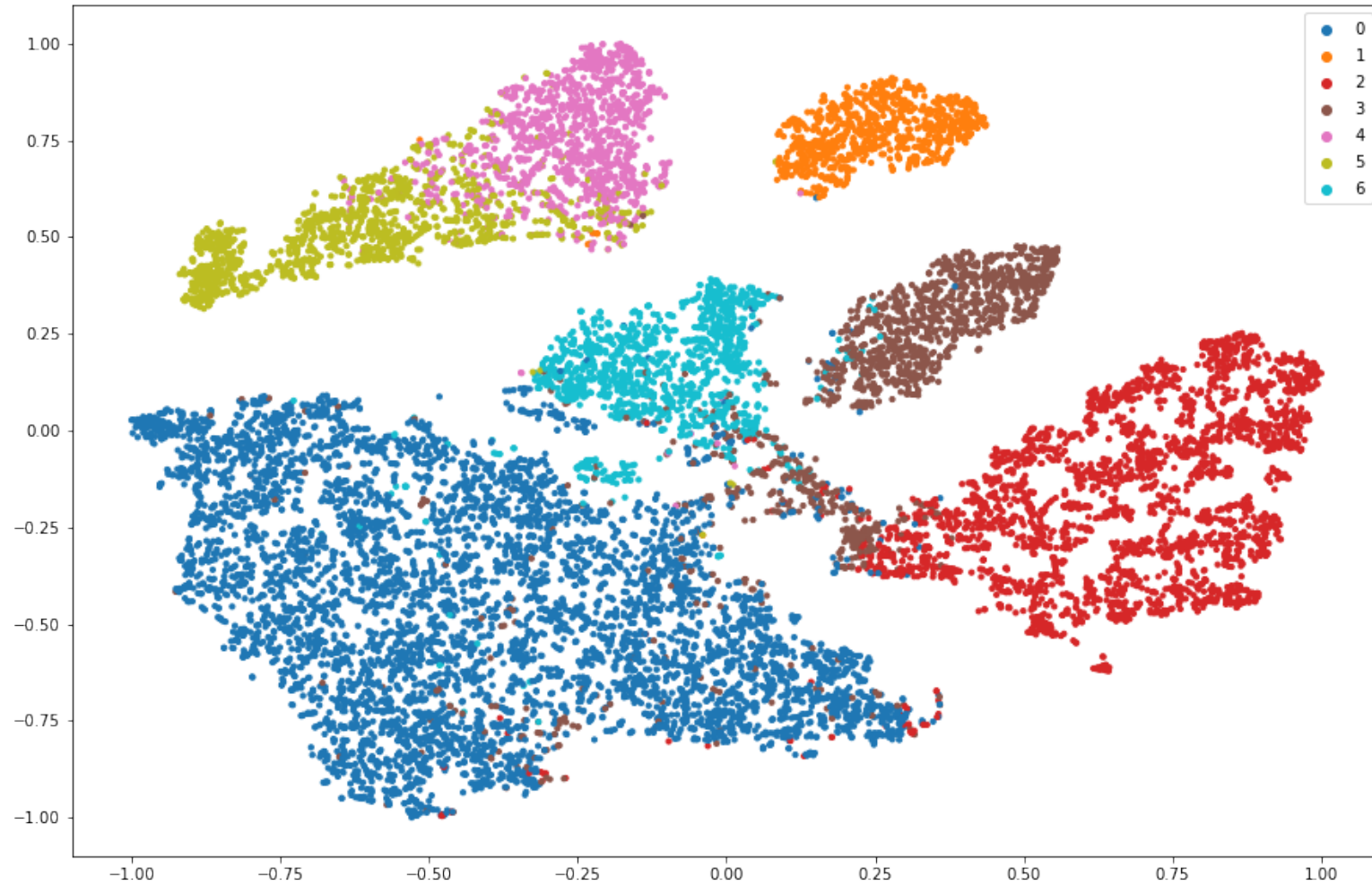


Variational AutoEncoder – MNIST

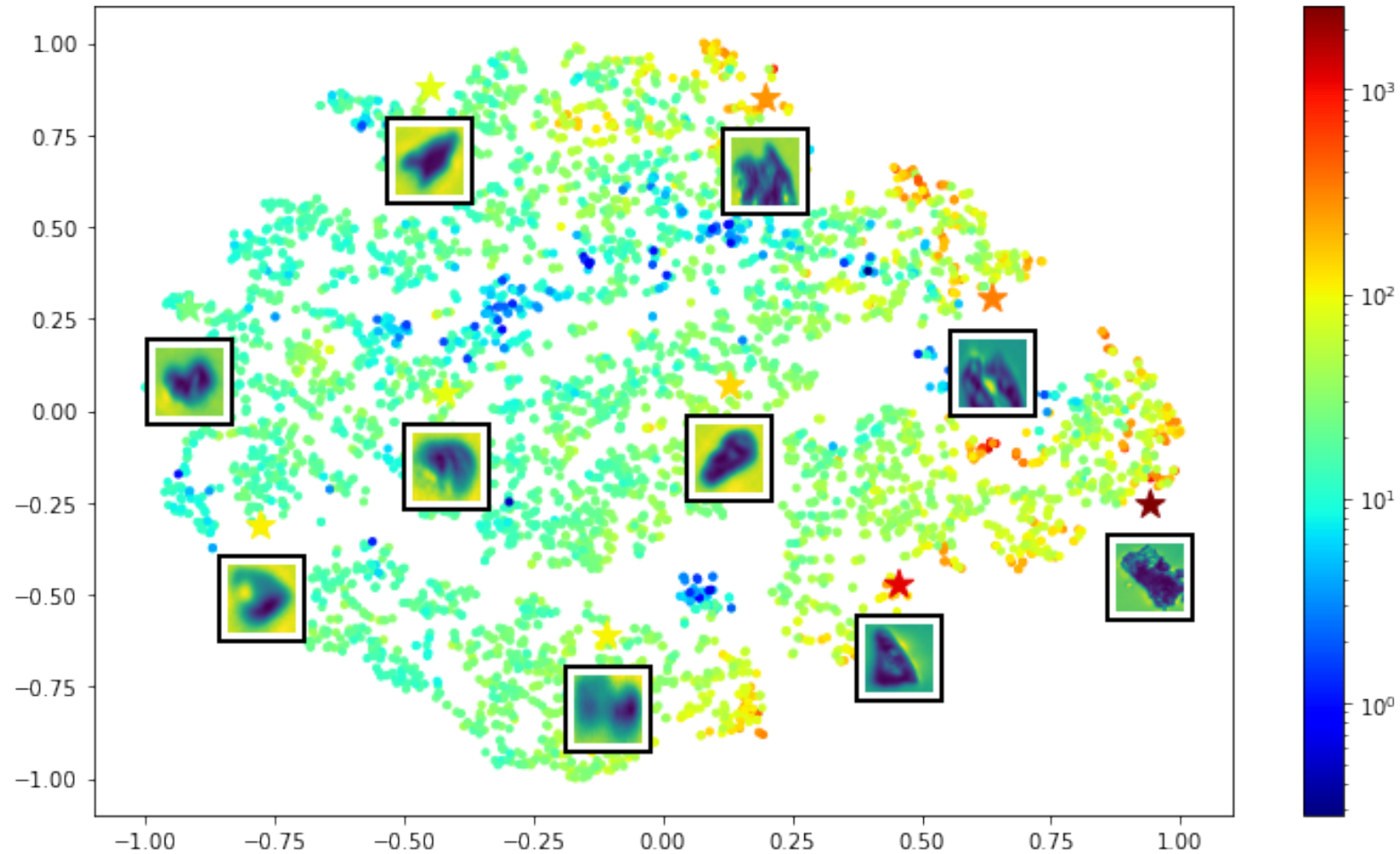
Mapping the latent space of a Variational AutoEncoder trained on the MNIST Train Dataset



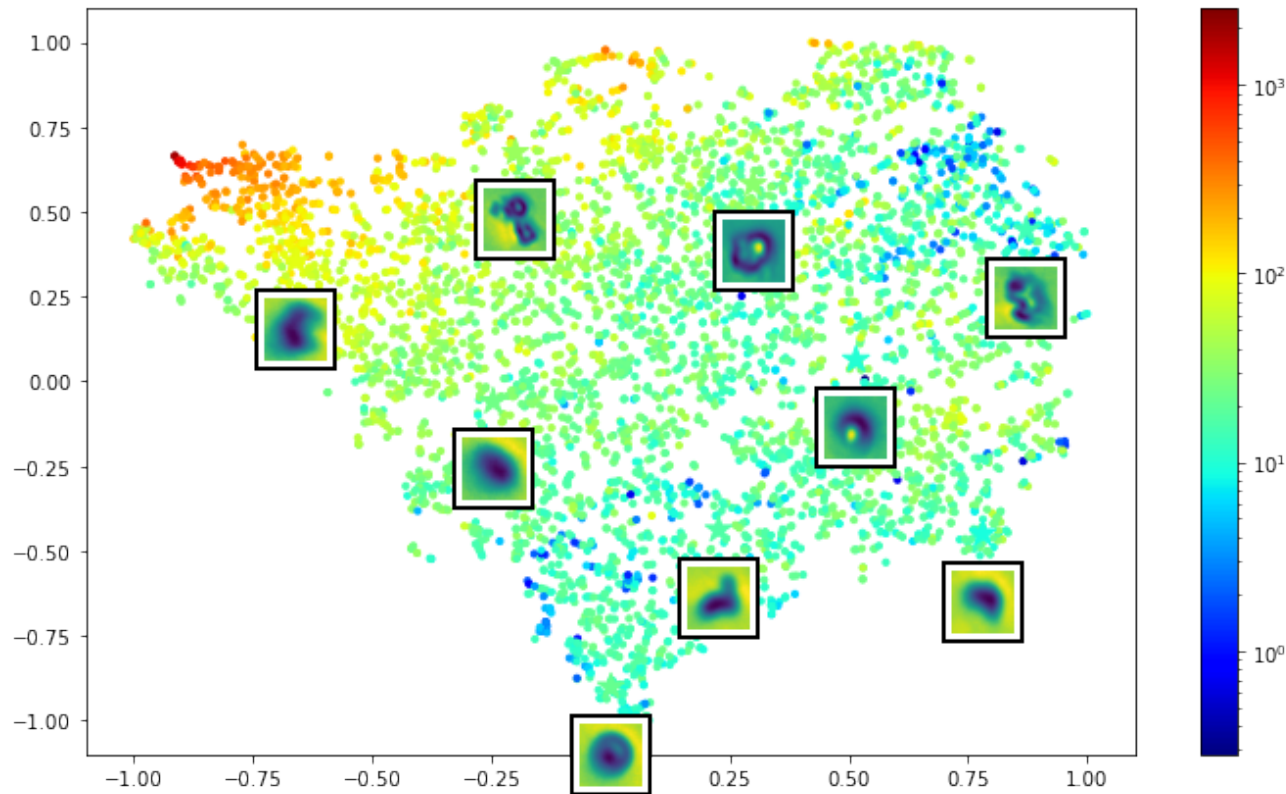
CNN-Classifier on Artificial Dataset



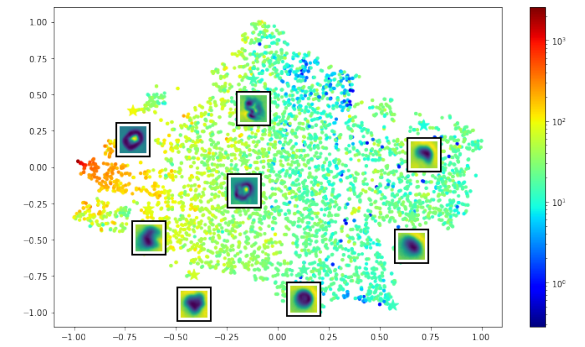
CNN-Classifier on Peruvian Dataset



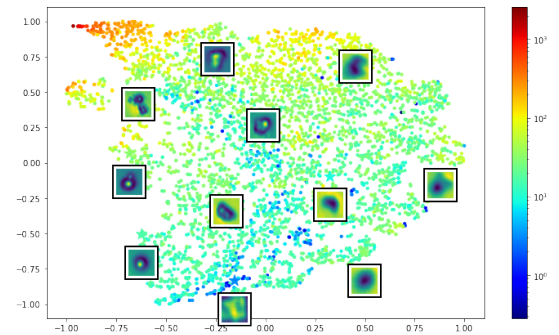
AutoEncoders on Peruvian Data



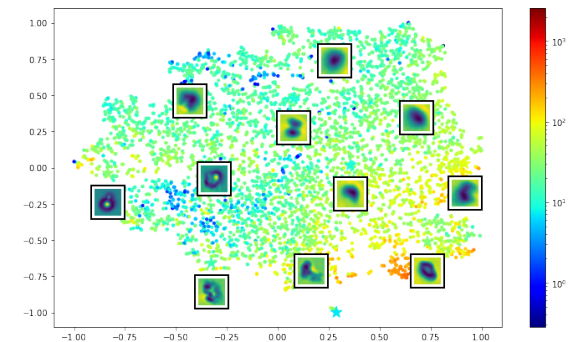
AutoEncoder with mean correction



AutoEncoder plain

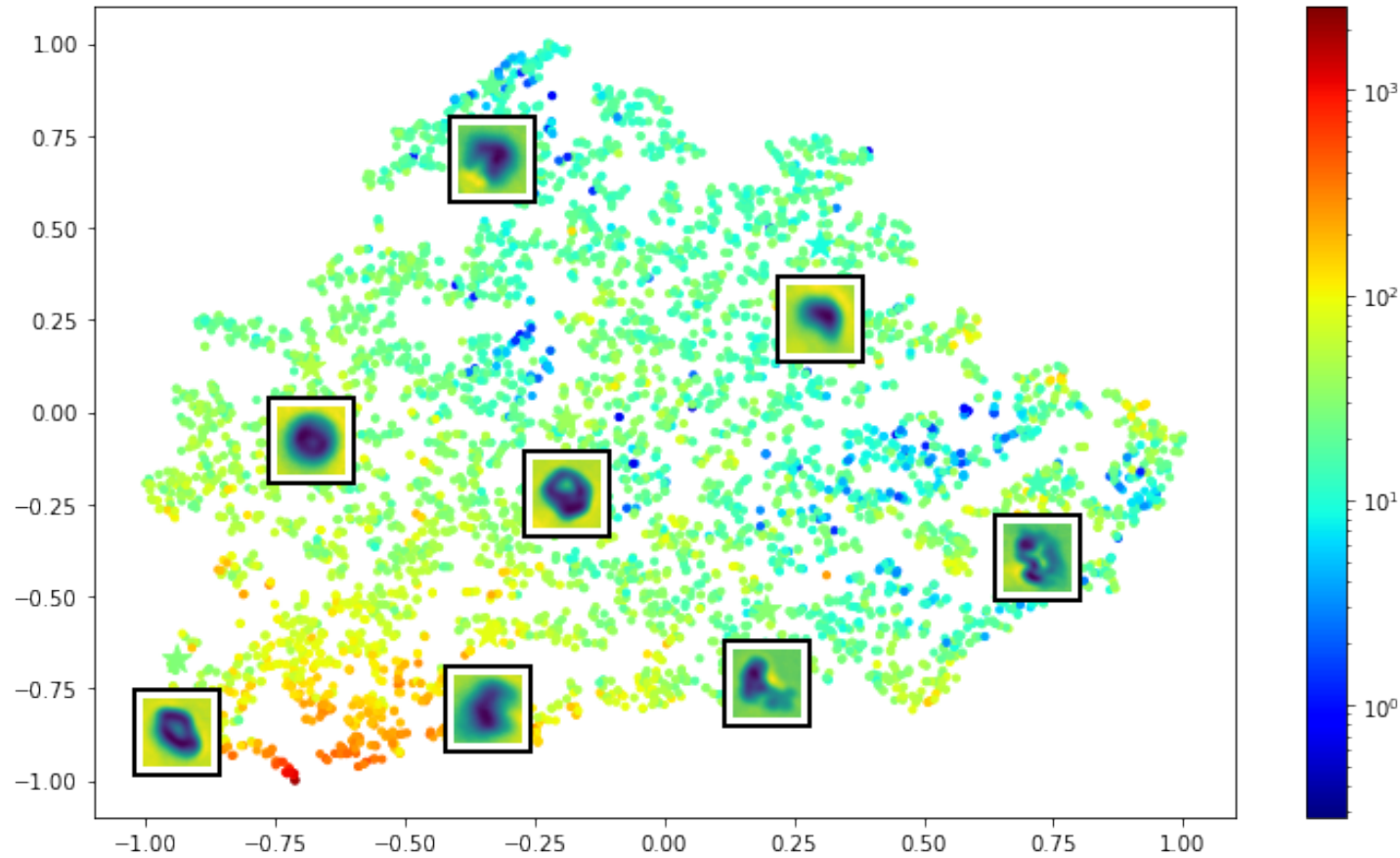


AutoEncoder with mean correction and no Dust

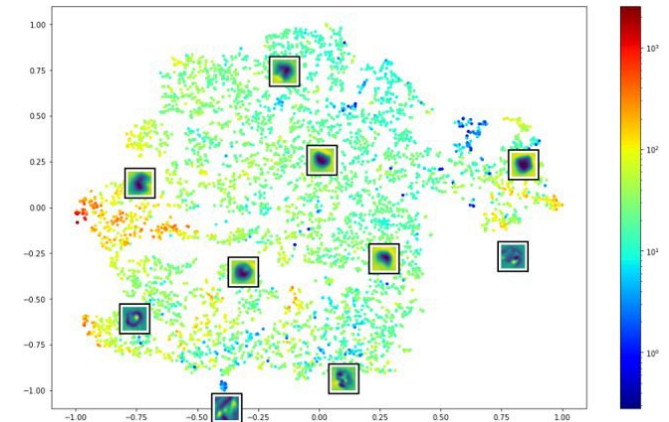


AutoEncoder plain latent space 128

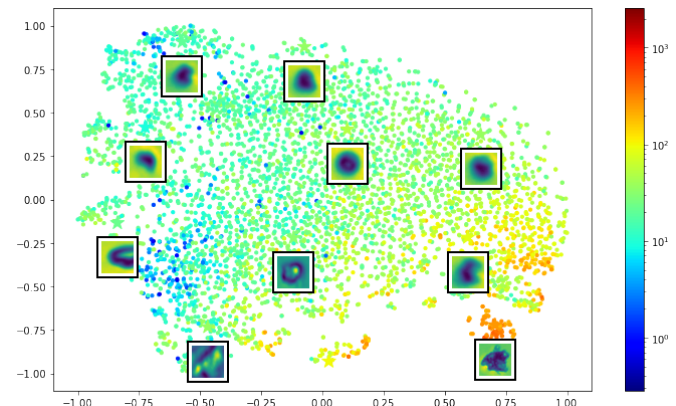
Variational AutoEncoders on Peruvian Dataset



Variational AutoEncoder with mean correction



Planar flow model trained without artificial data



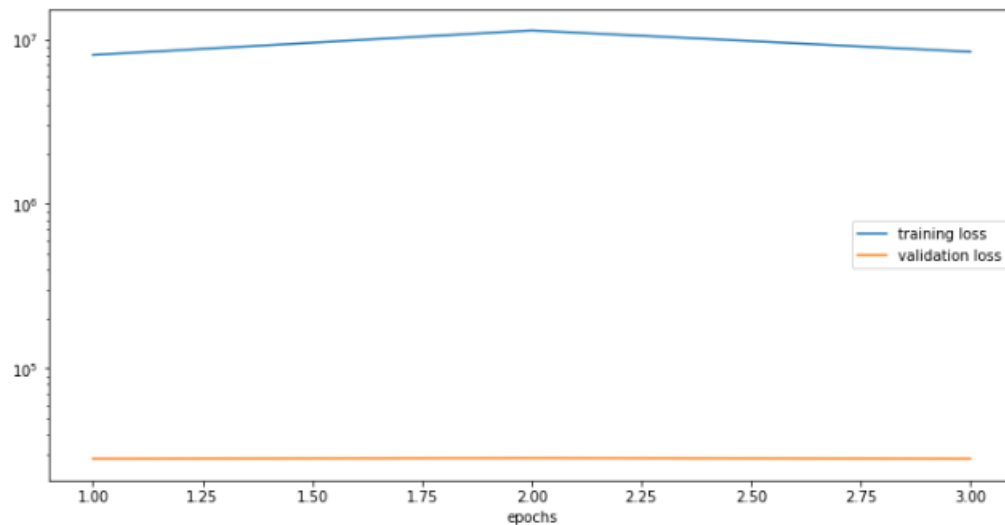
Planar flow trained with artificial data

Difficulties during training

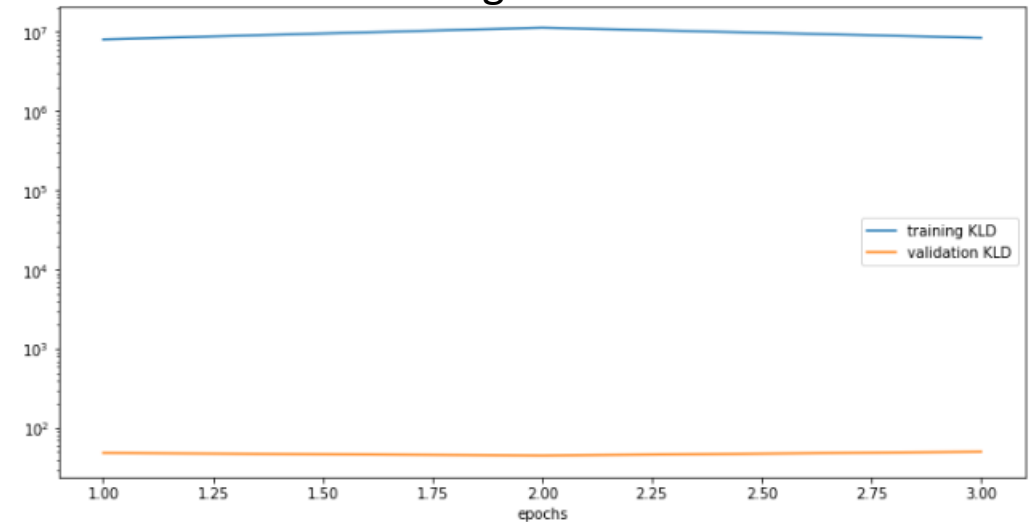
– Overfitting on initial parameters?

- KL divergence for some events in the Peru set explodes before training has even started and never seem to drop. Most likely overfitting encoder from classifier on the artificial set. And yes! Problem solved by reducing epochs on pretrained encoder.

Total loss

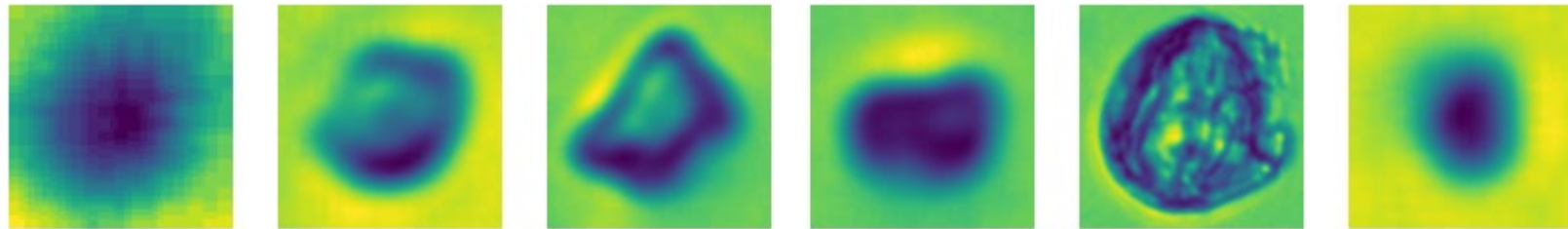


Kullback Liebler divergence

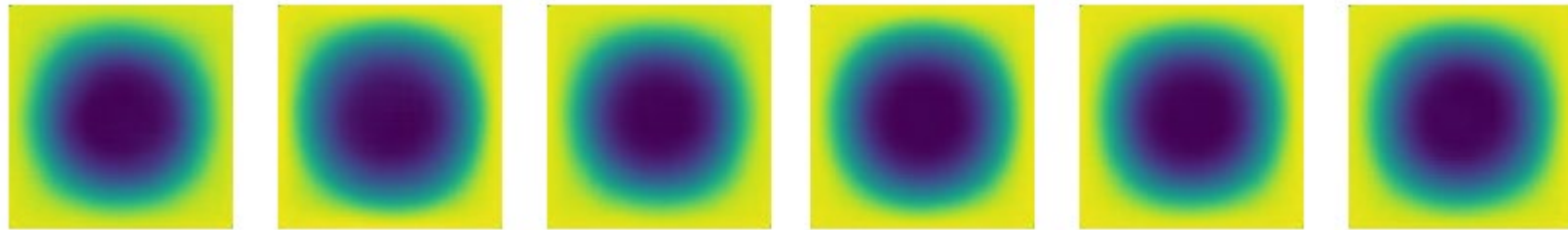


Difficulties during training – mode collapse?

Input images



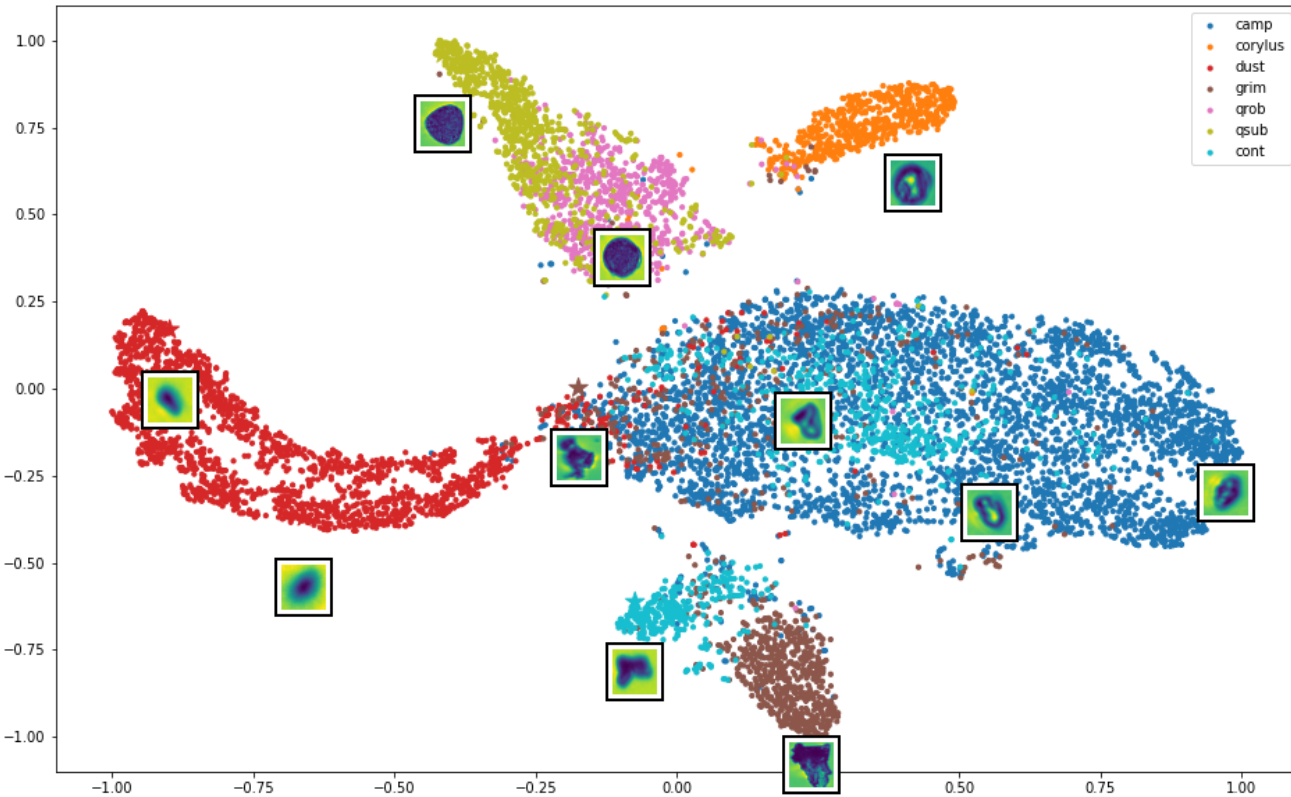
Output images



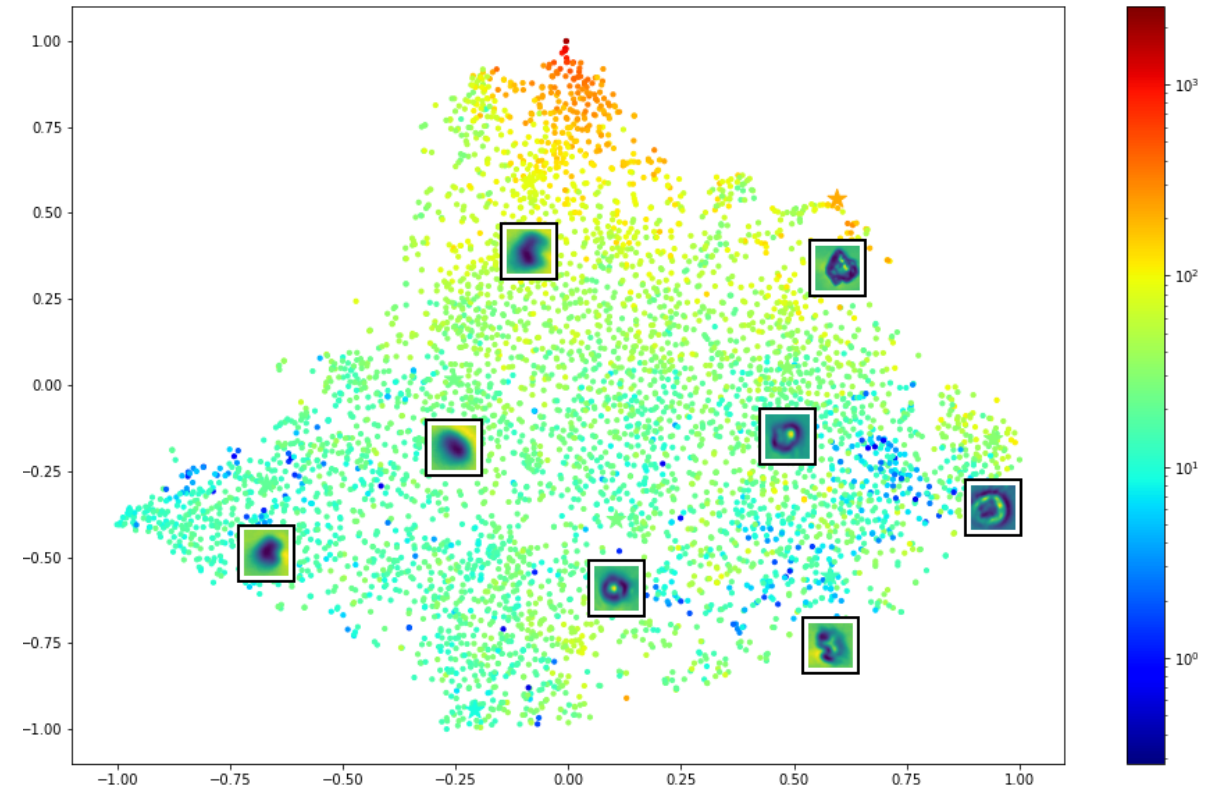
KLD was about 2 vs about 40 in "good" trainings indicating information depleted latent space.

Difficulties with overfitting

Artificial dataset

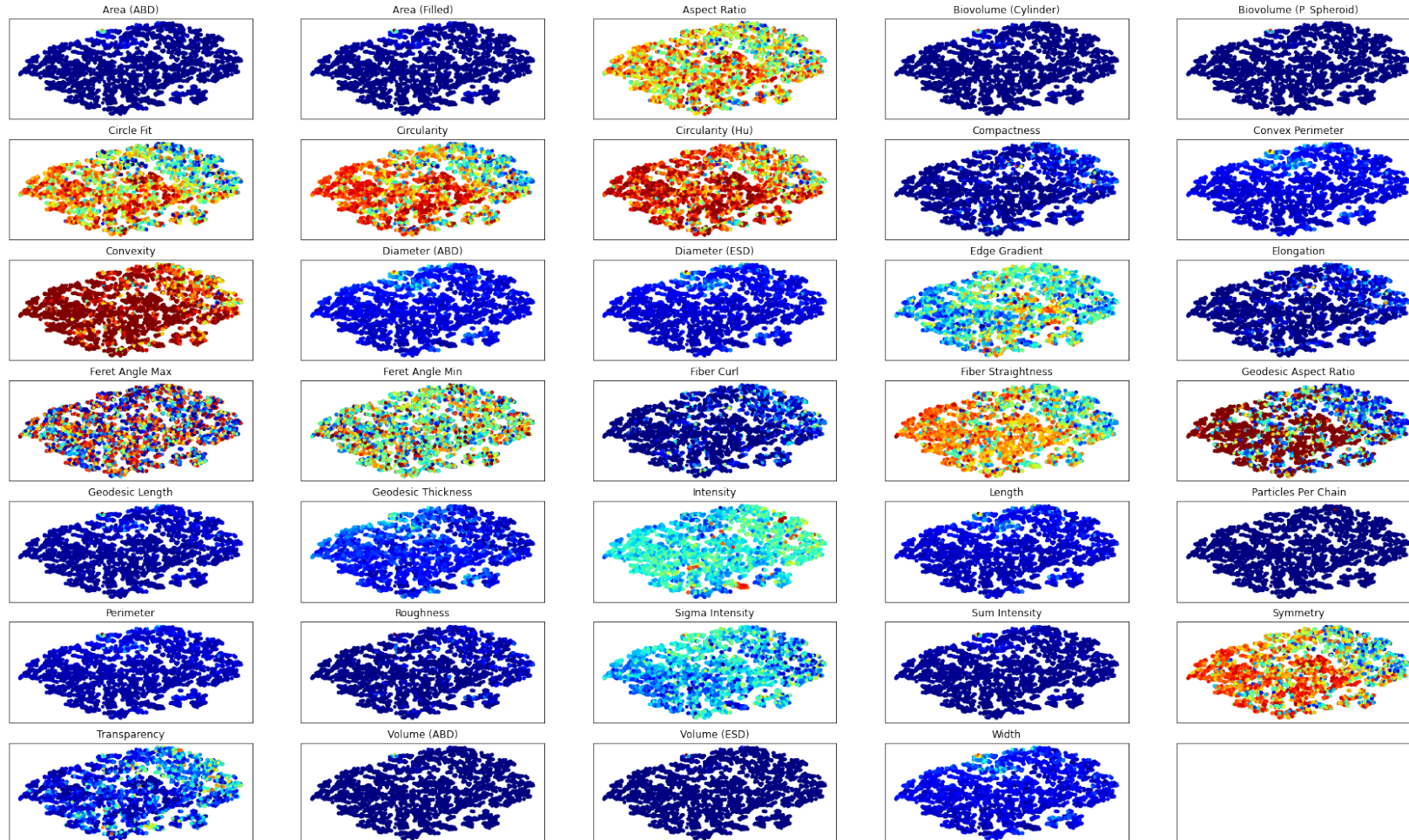


Peruvian set



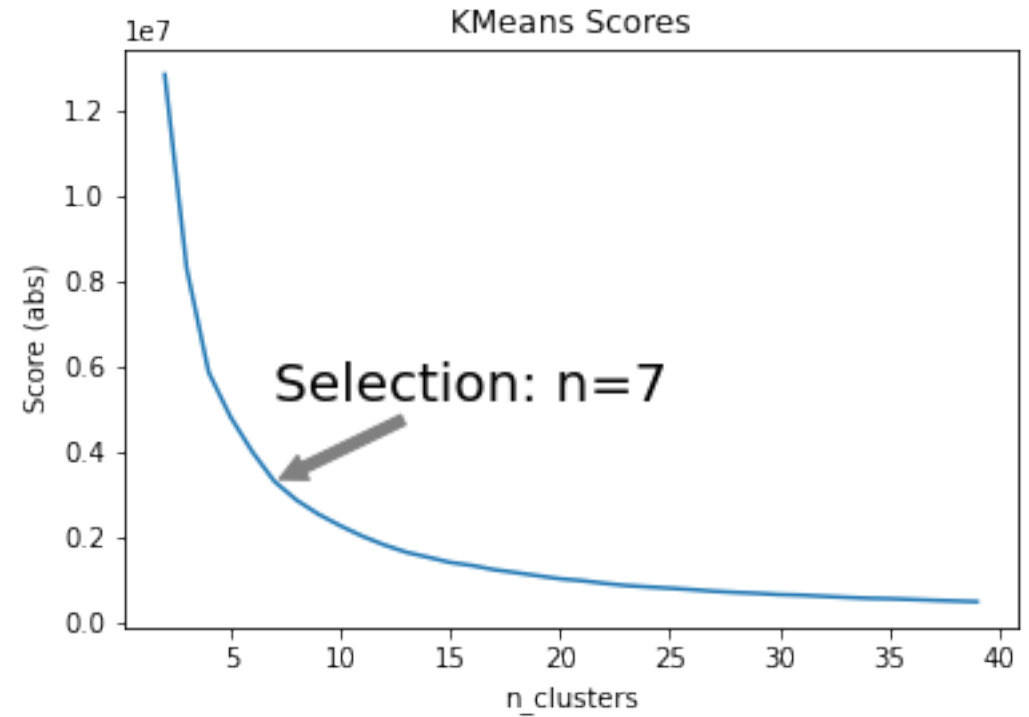
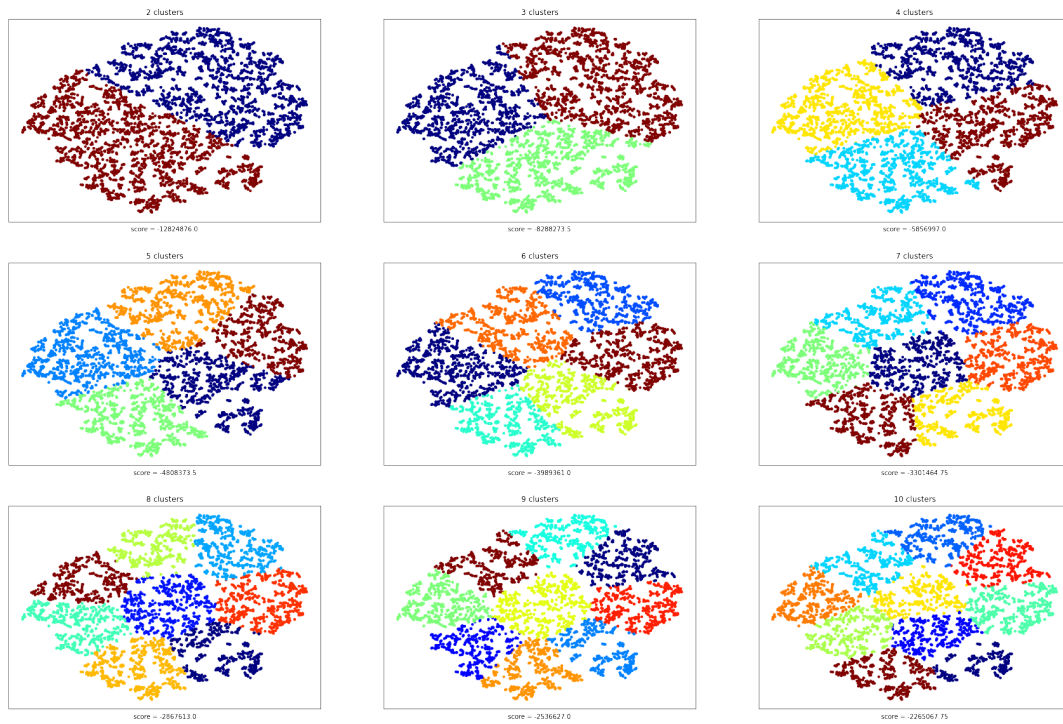
Fitted too much on labeled data?

Plotting scalars

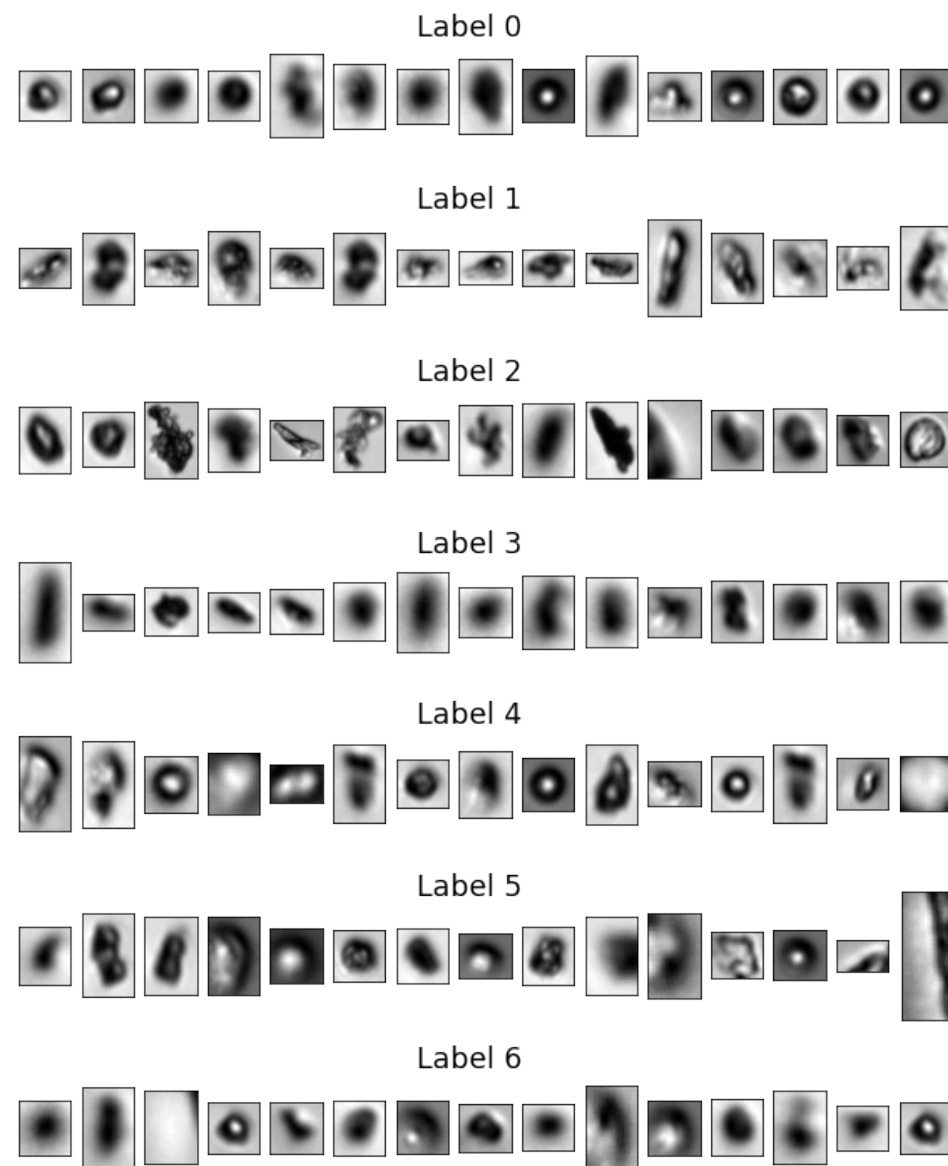
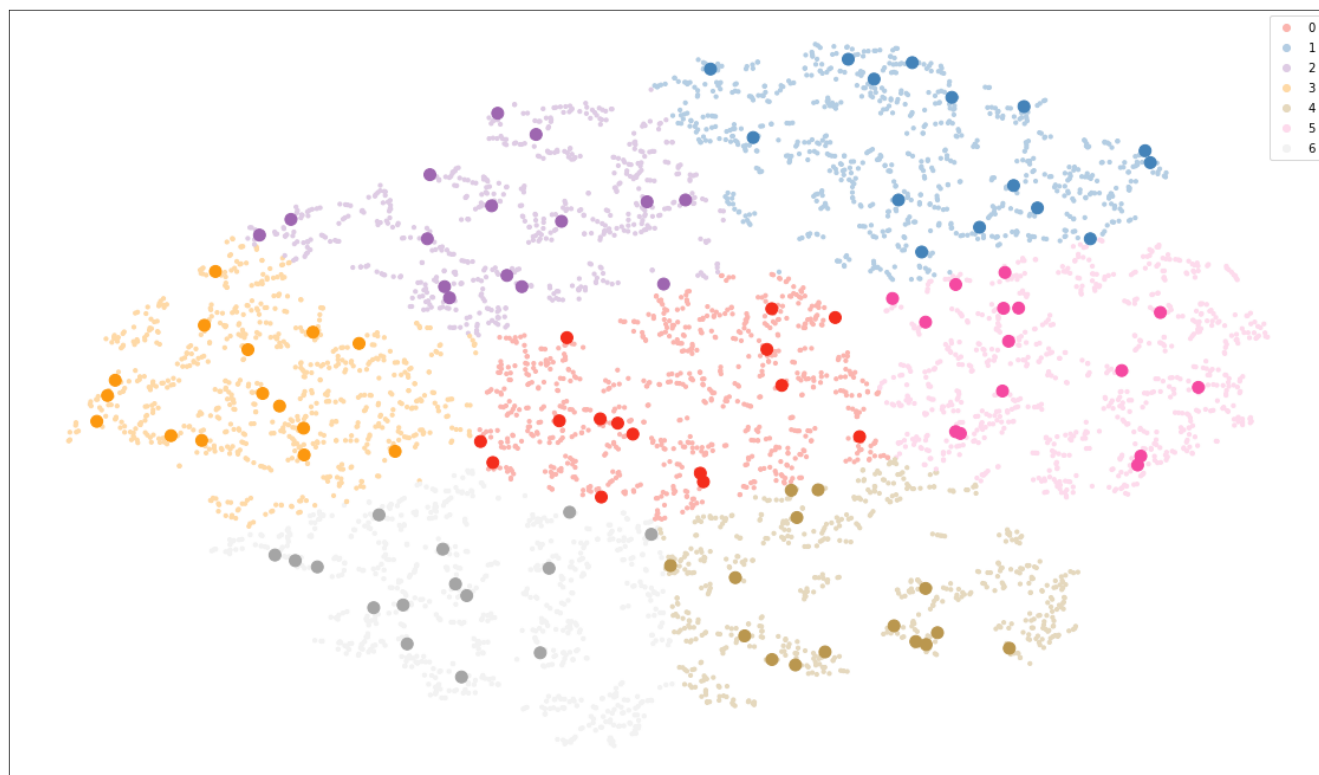


KMeans Clustering

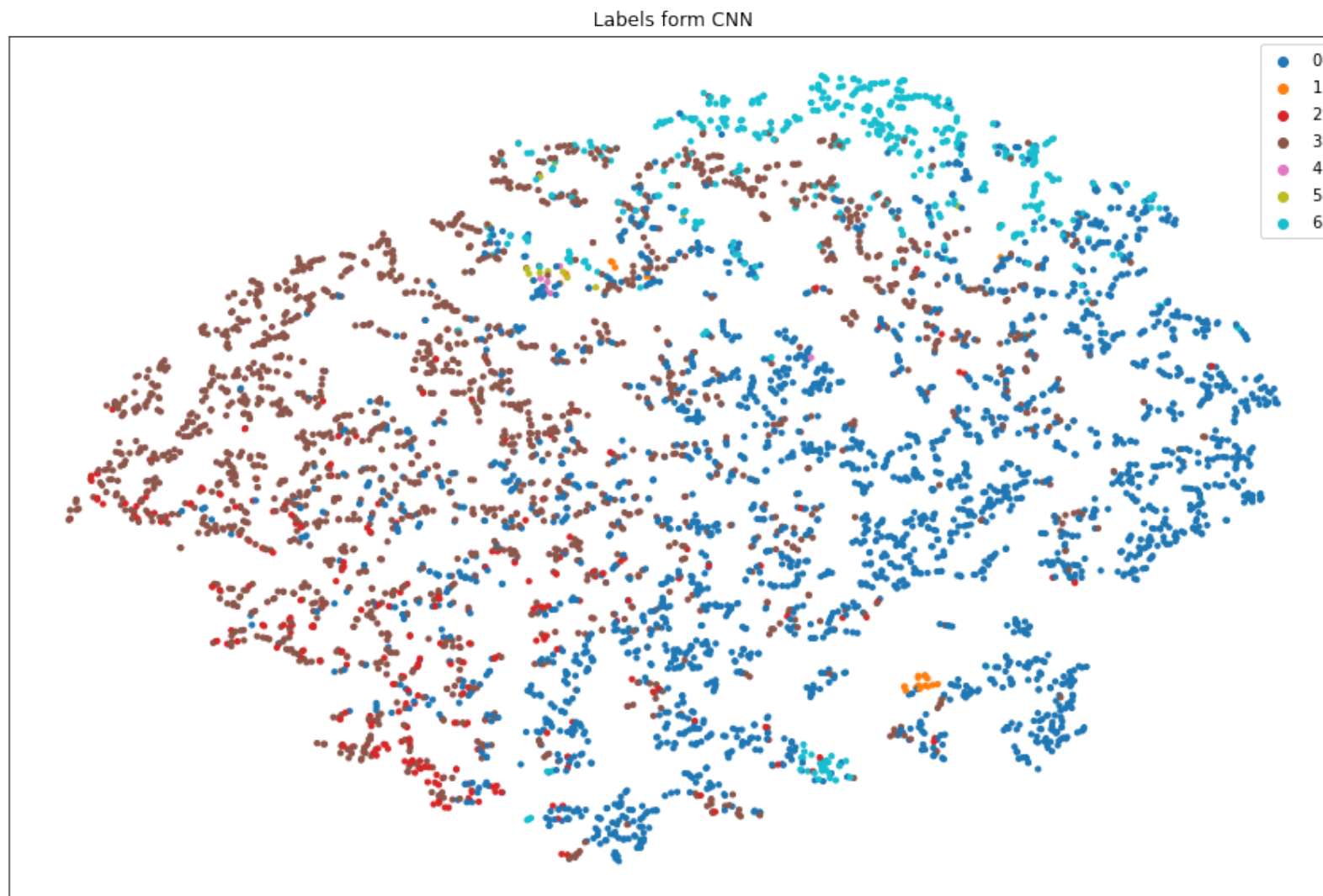
KMeans clustering



KMeans for n=7

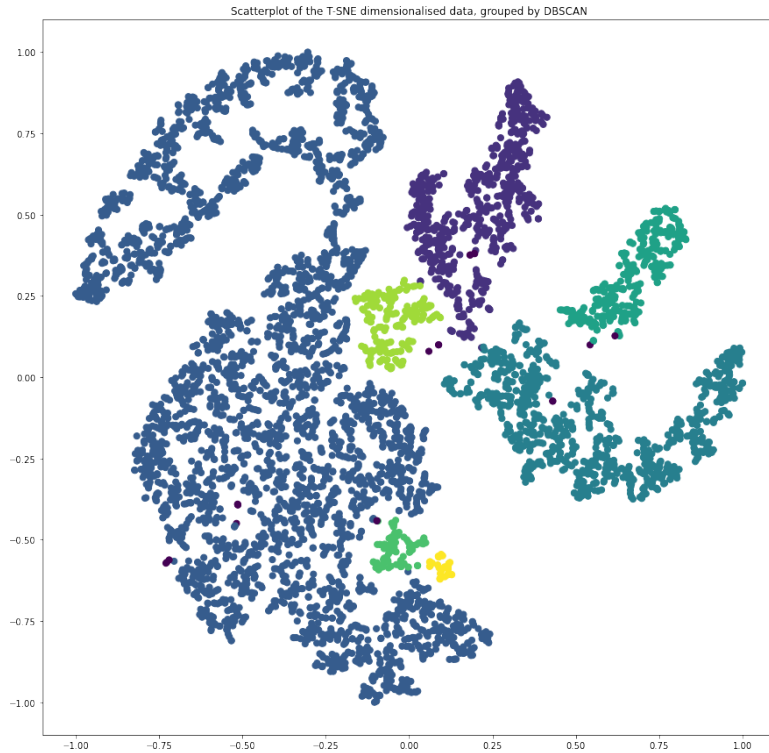


Overlaying labels from CNN-Classifier

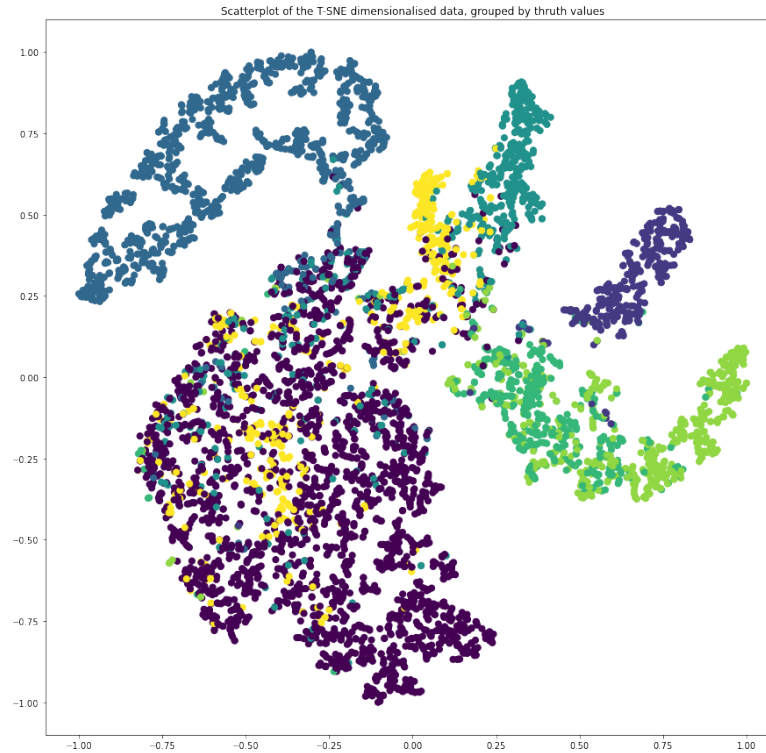


DBSCAN

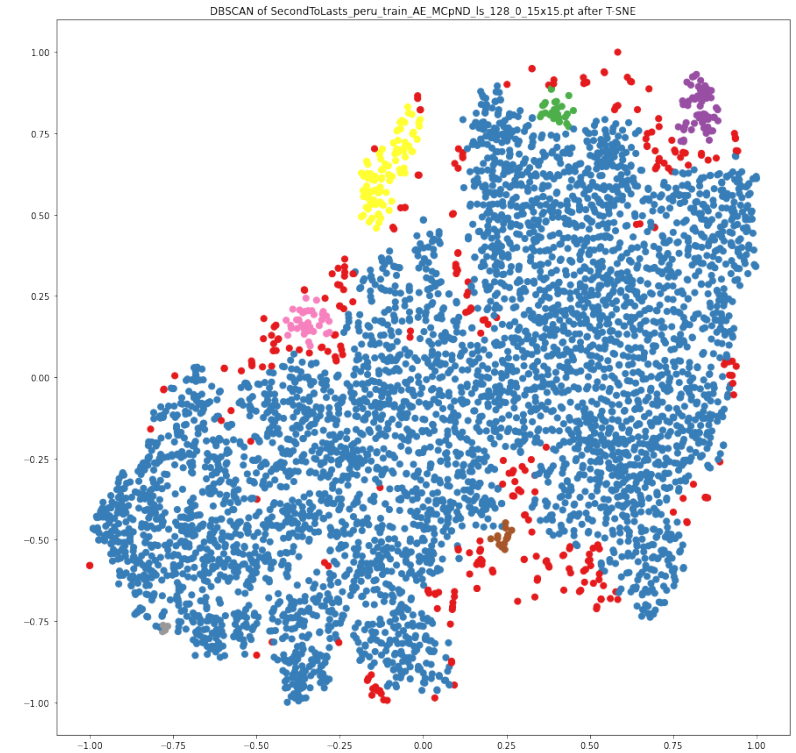
- DBSCAN on the AE/VAE tensors allows for partial clustering



Artificial Dataset (Prediction)

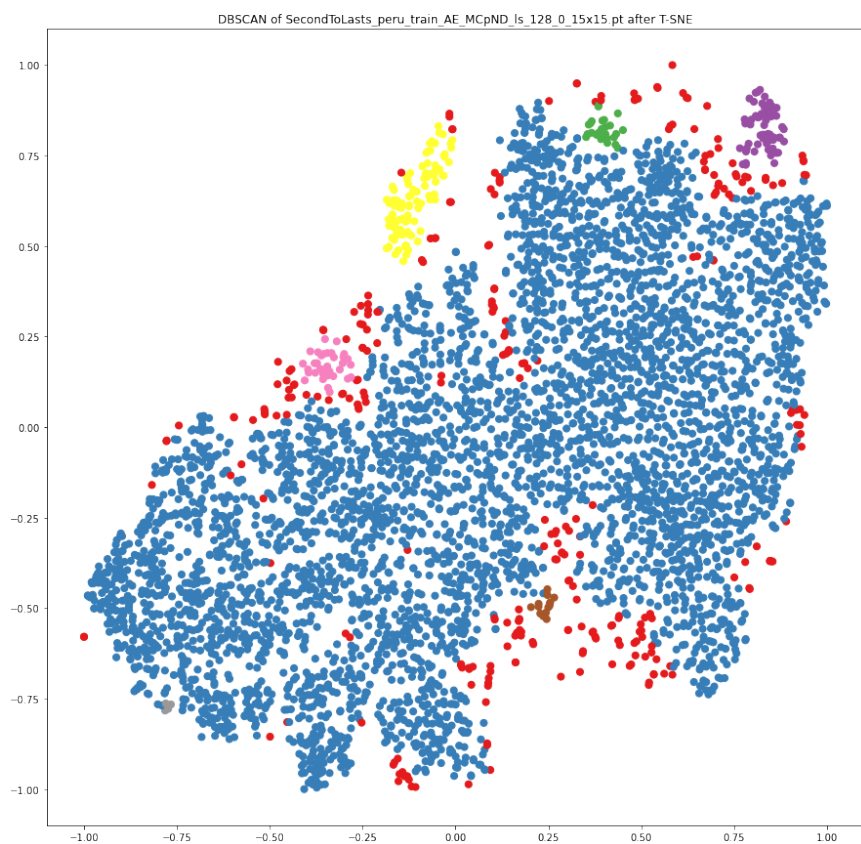


Artificial Dataset (Truth)

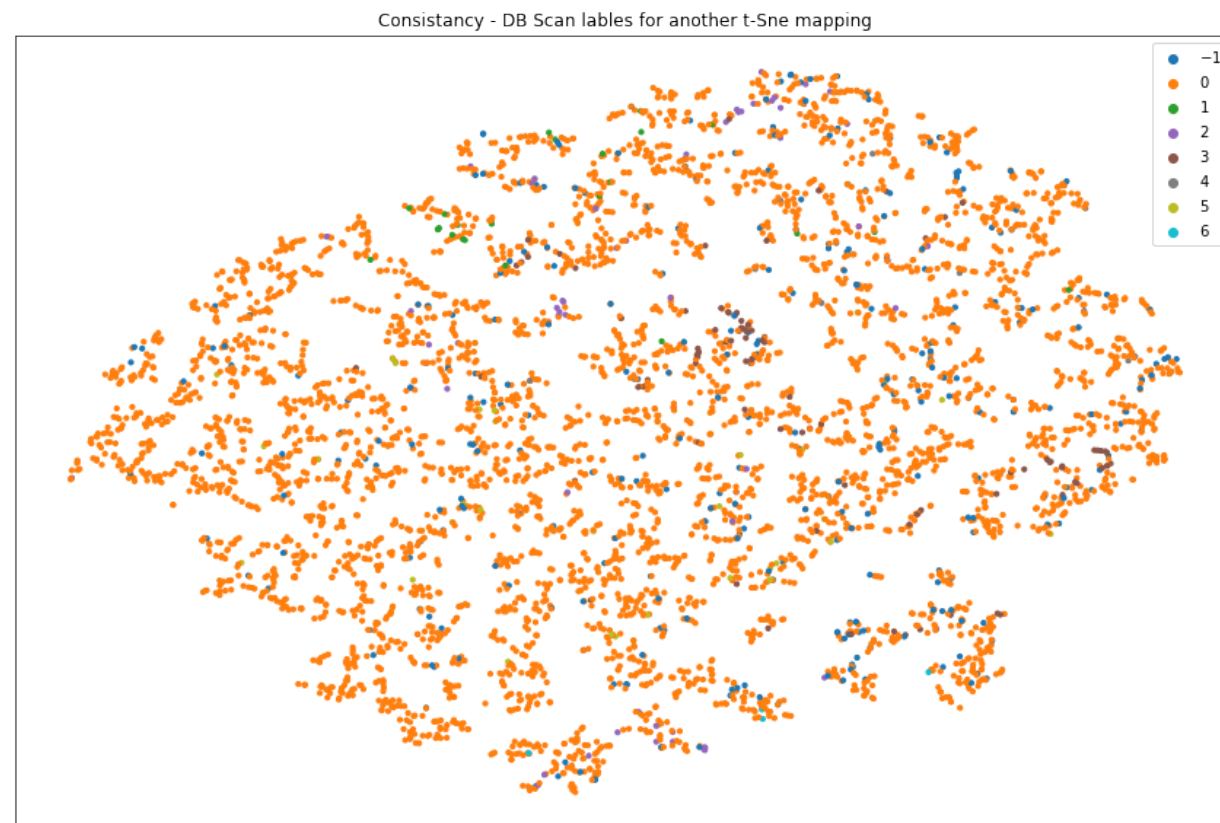


Peruvian Dataset

DBSCAN

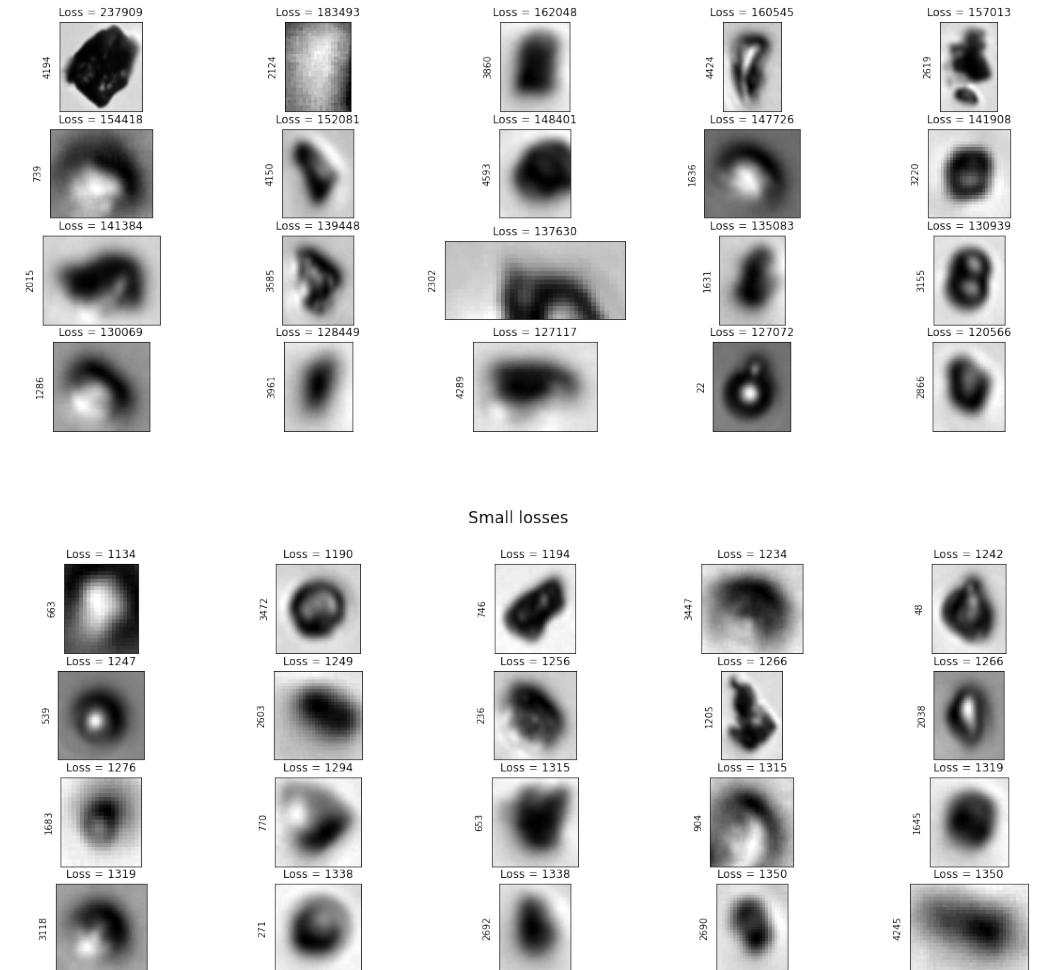
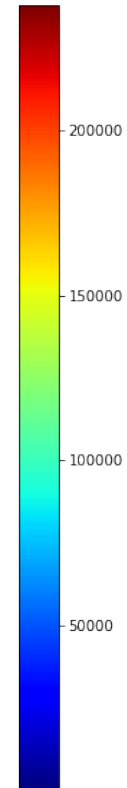
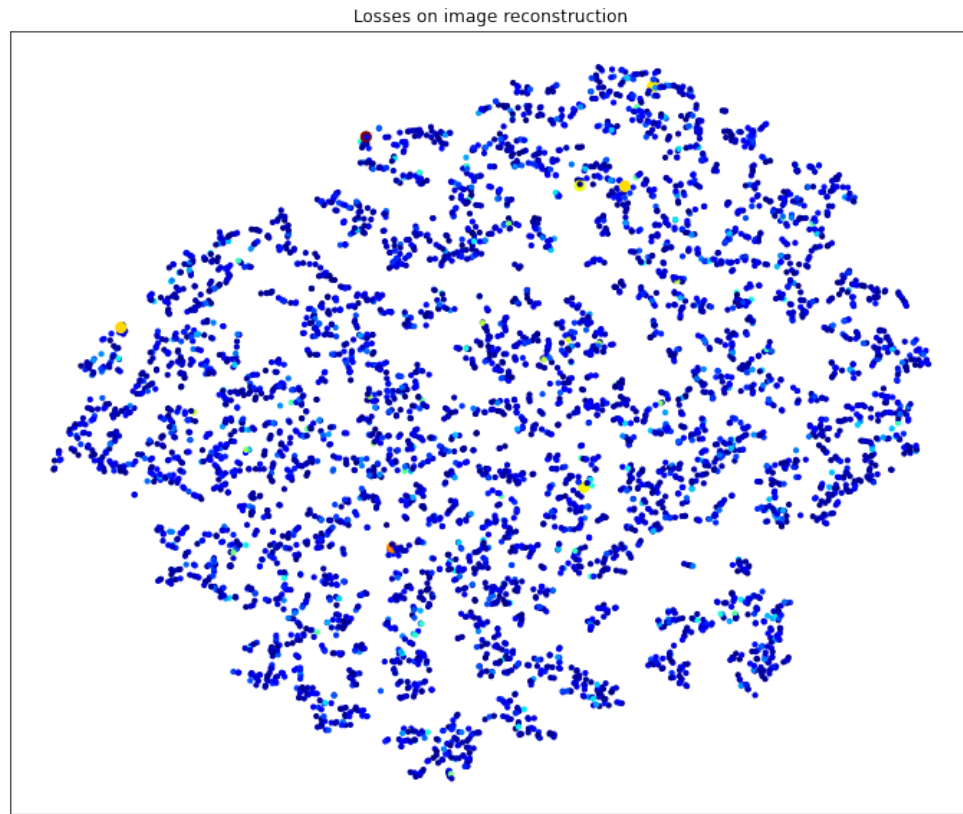


DBSCAN on the AE/VAE tensors allows for partial clustering

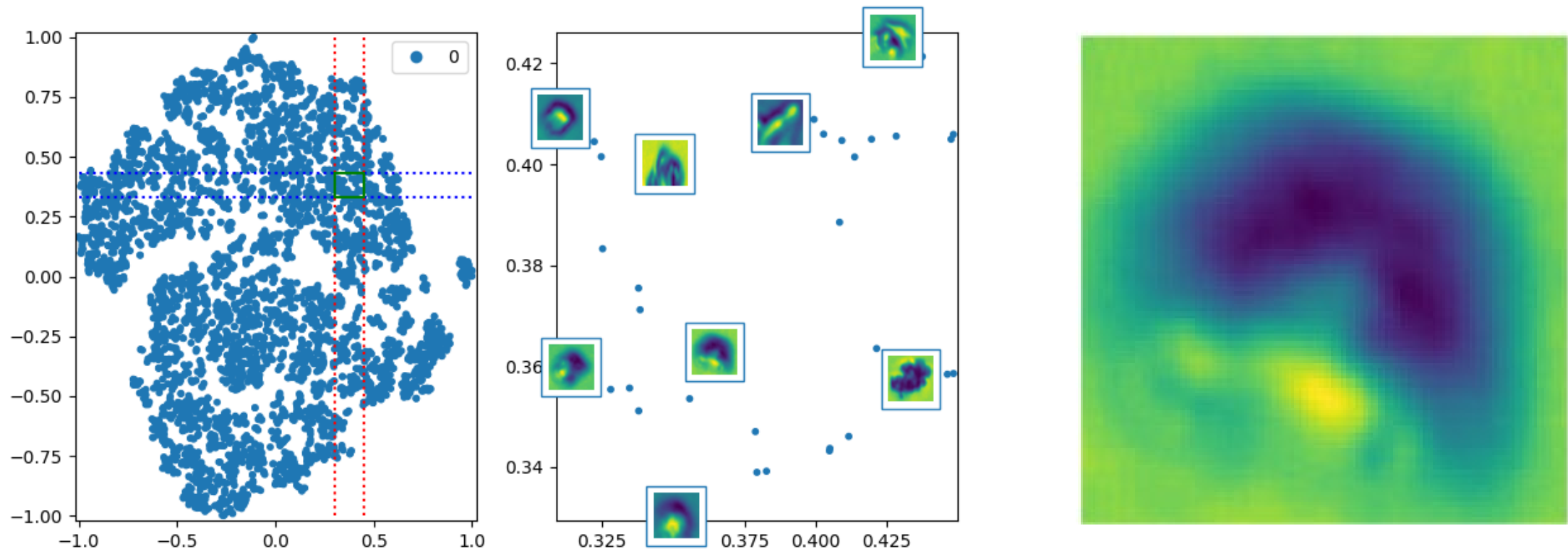


Consistency between mappings

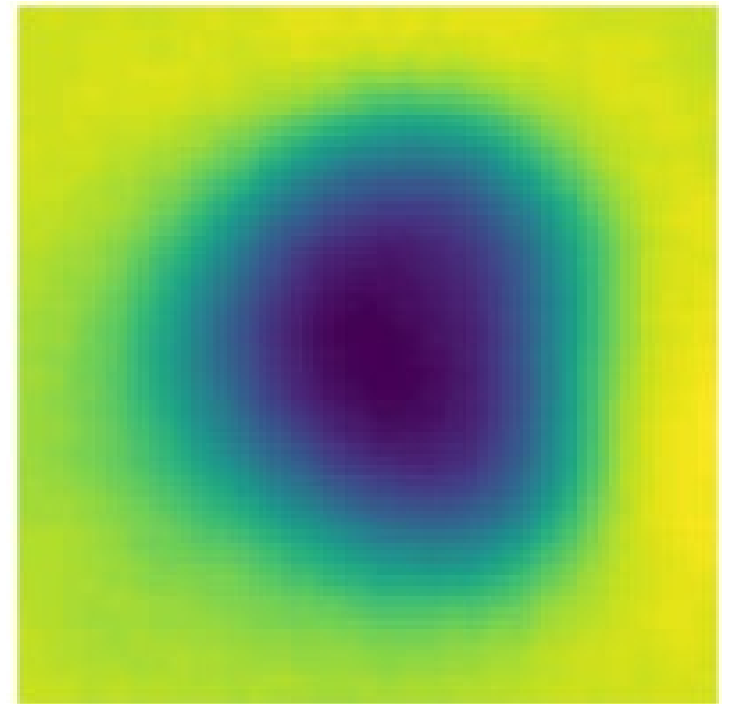
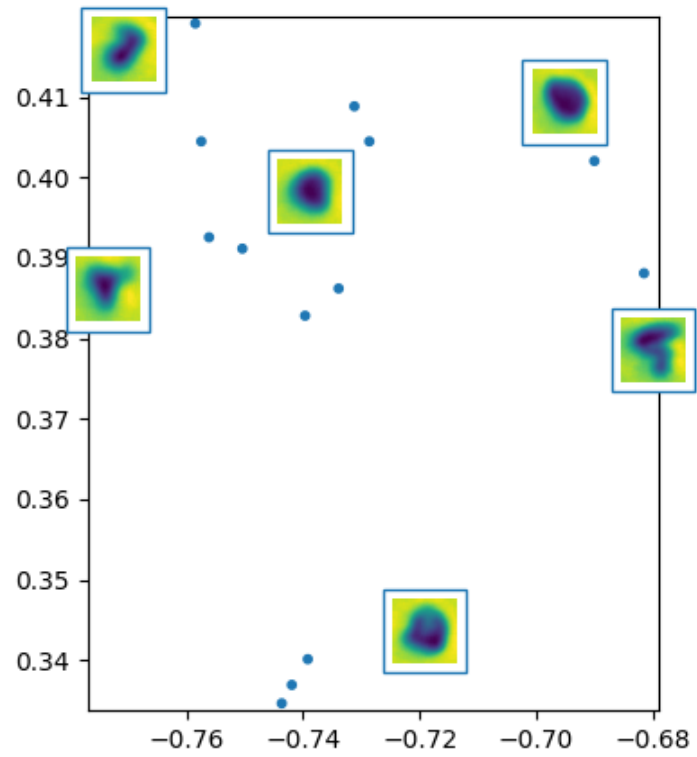
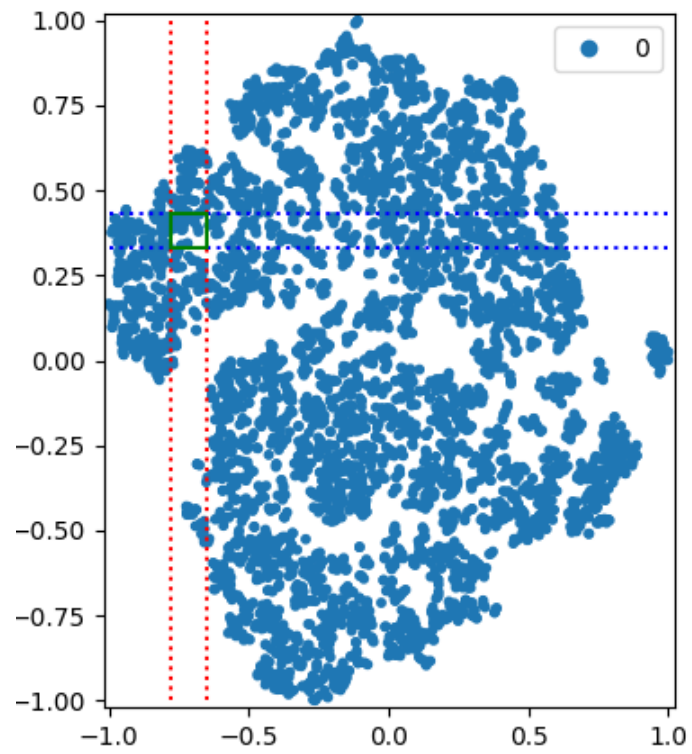
Plotting the loss



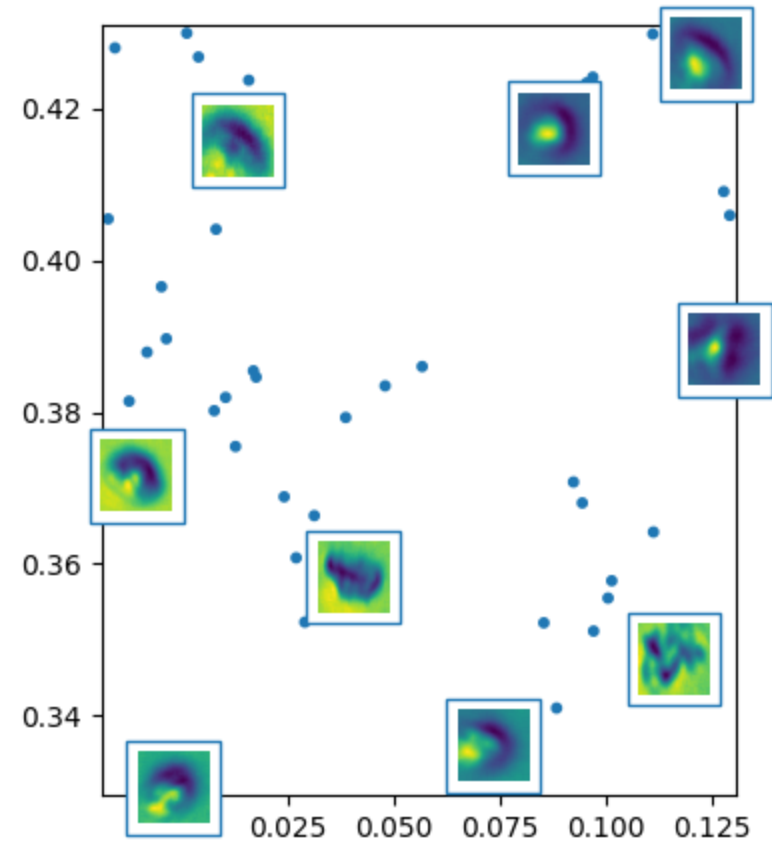
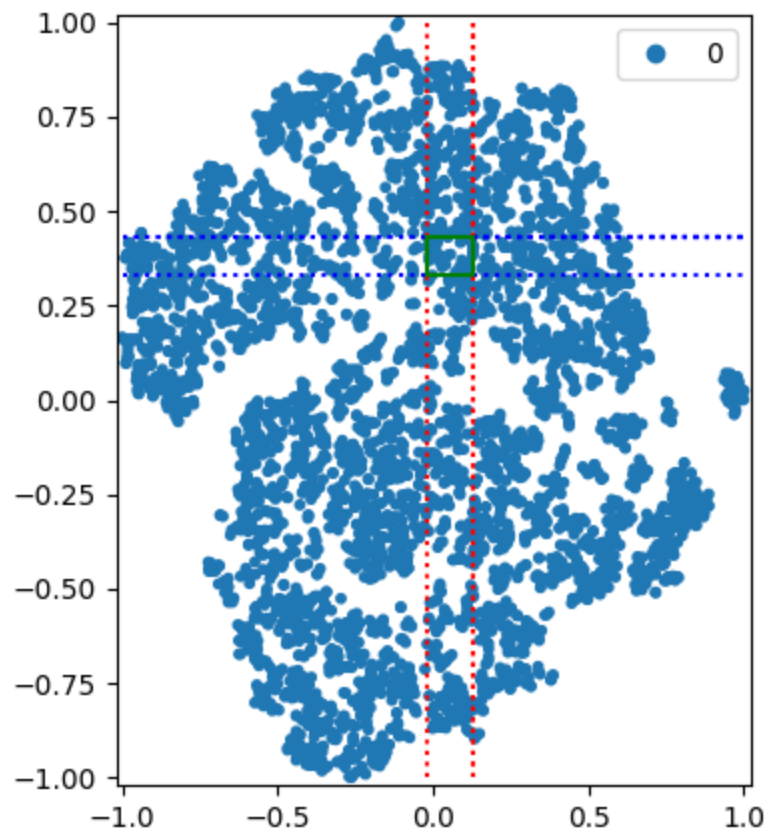
CNN classifier - Area 1



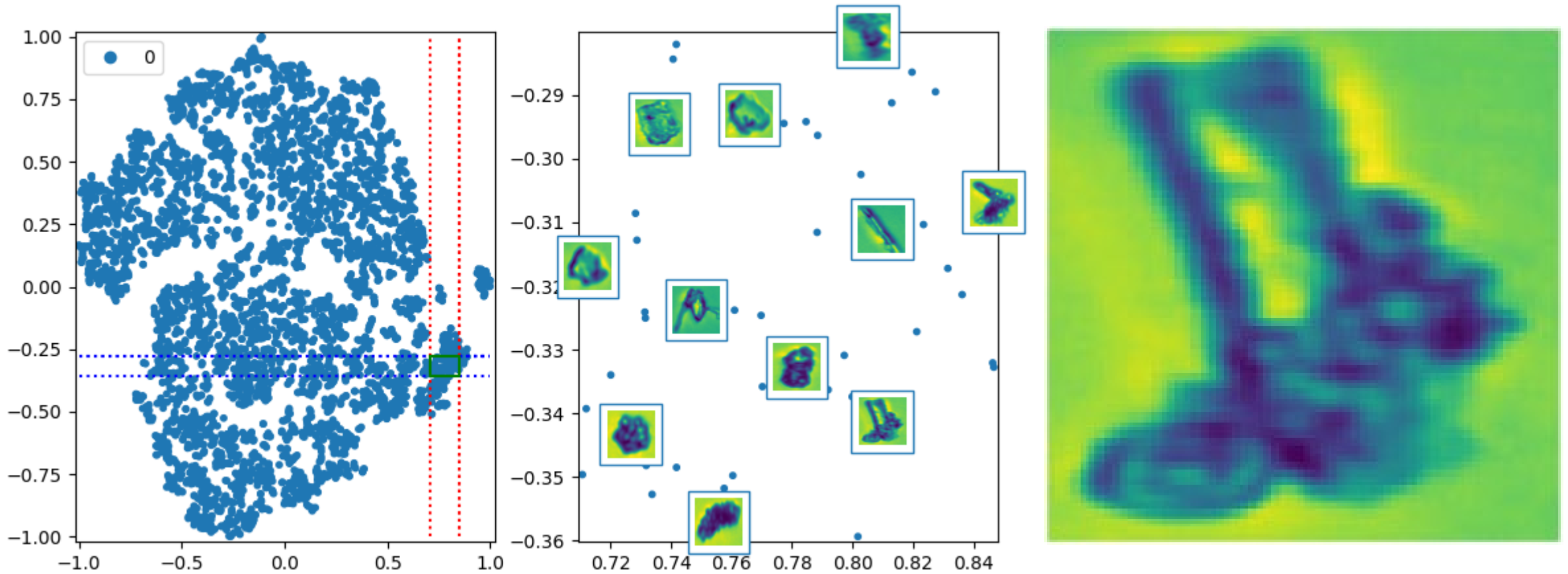
CNN classifier - Area 2



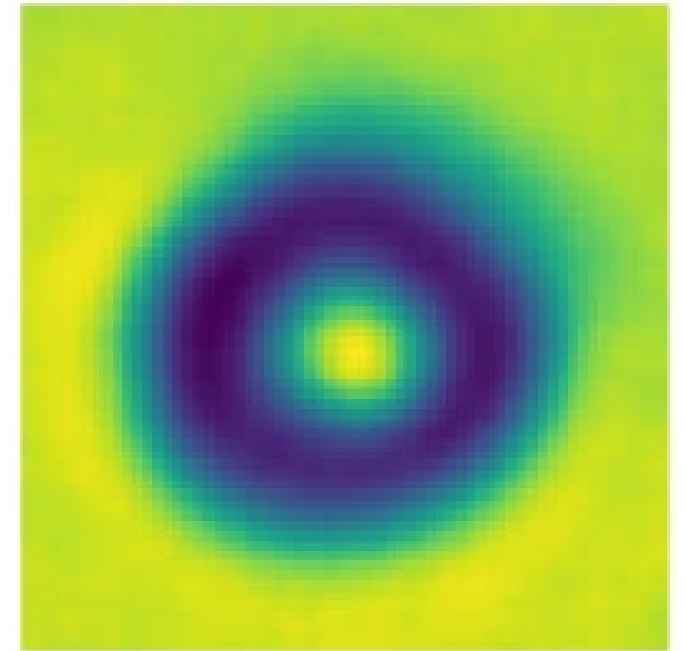
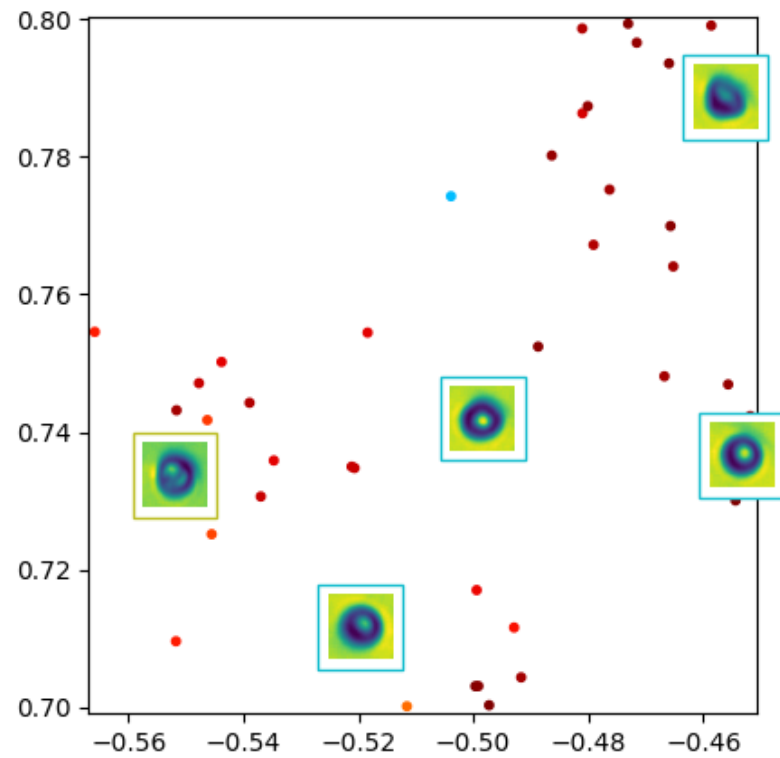
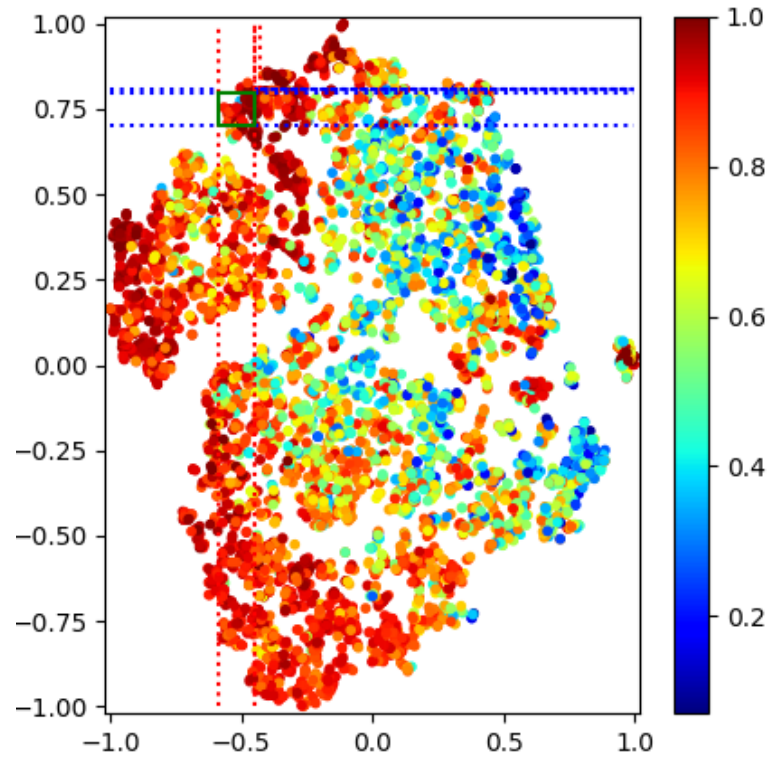
Between area 1 and 2



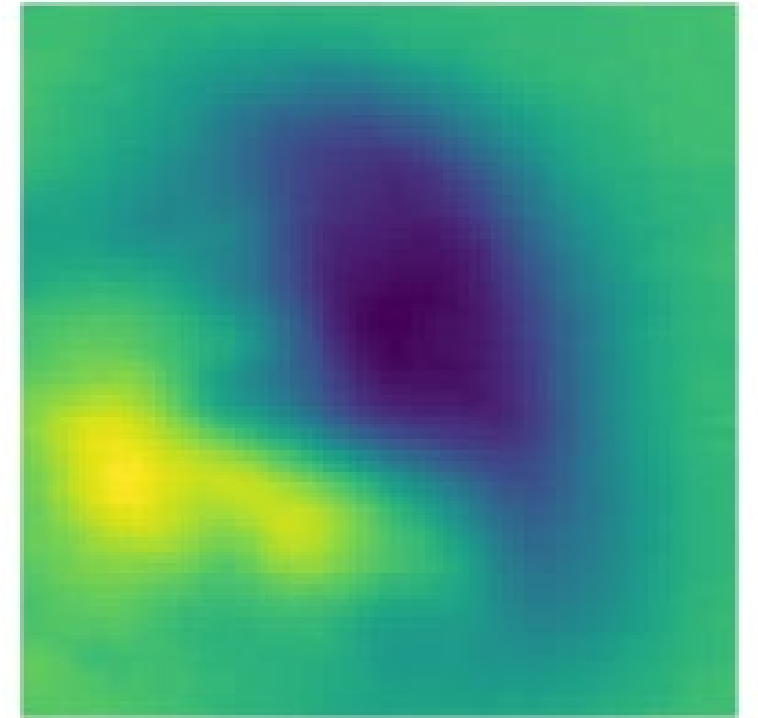
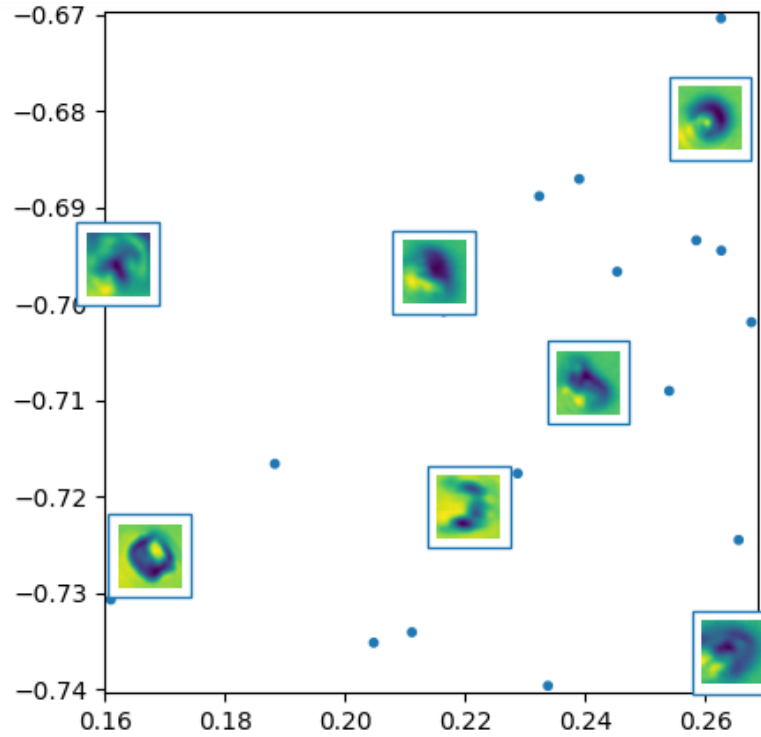
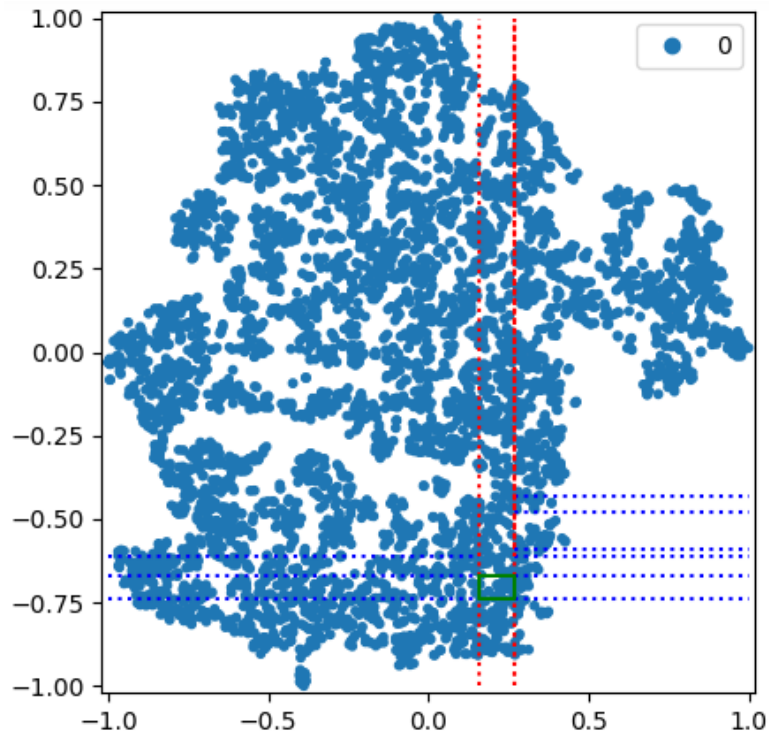
CNN classifier - Area 3



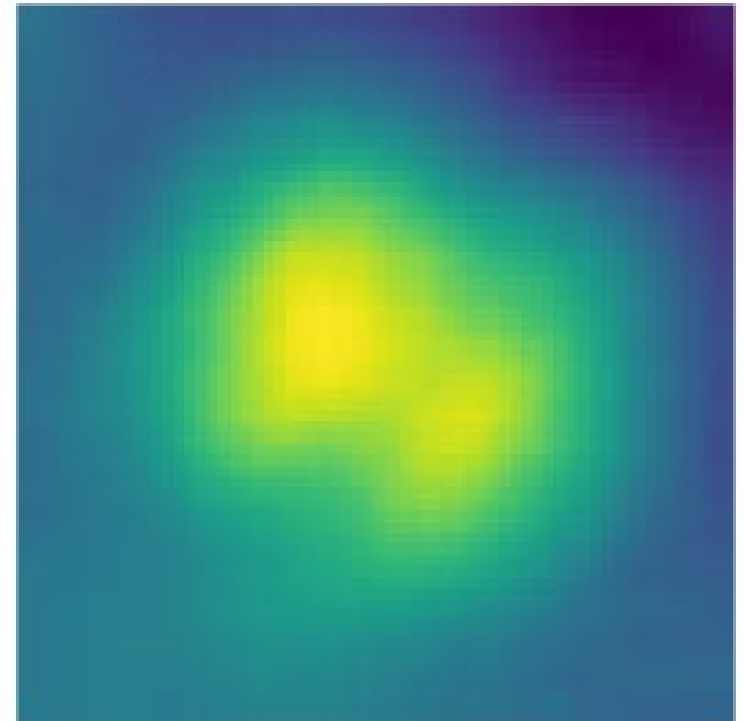
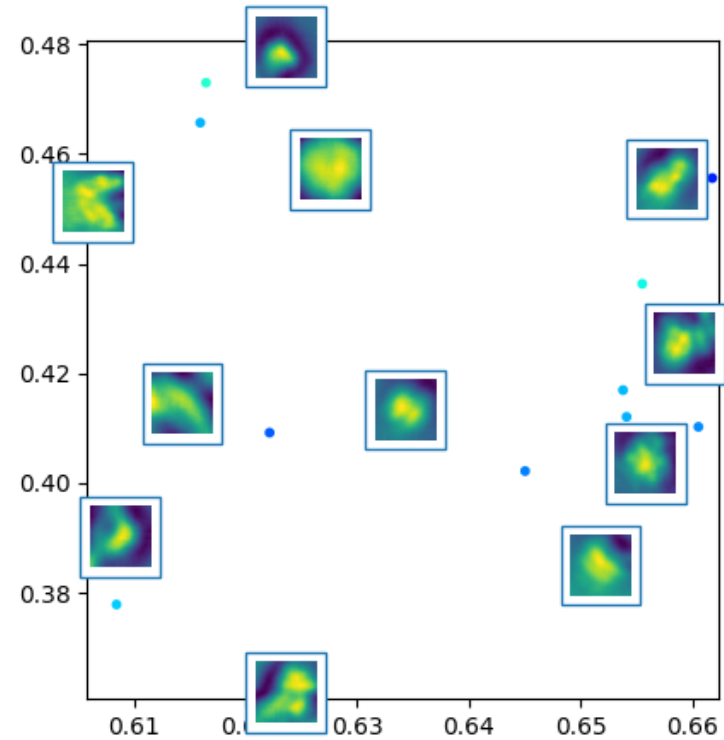
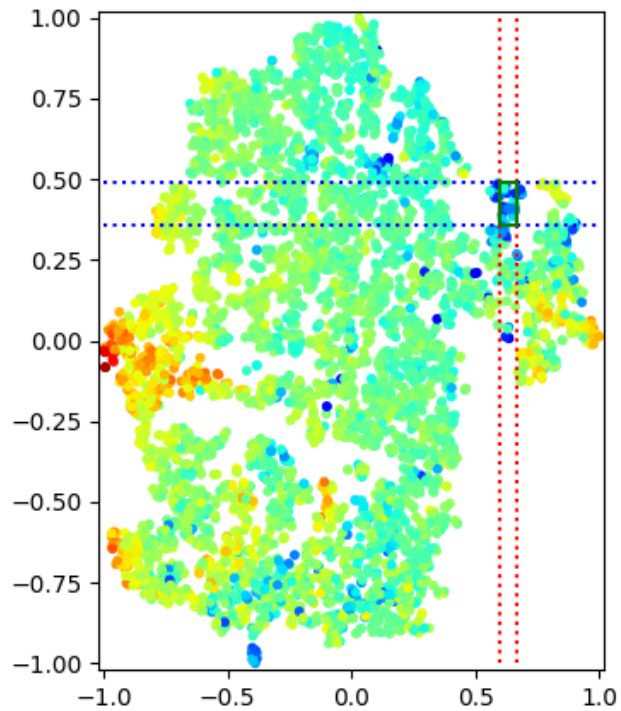
CNN classifier - Area 4



VAE, planar flow Area 1



VAE, planar flow area 2

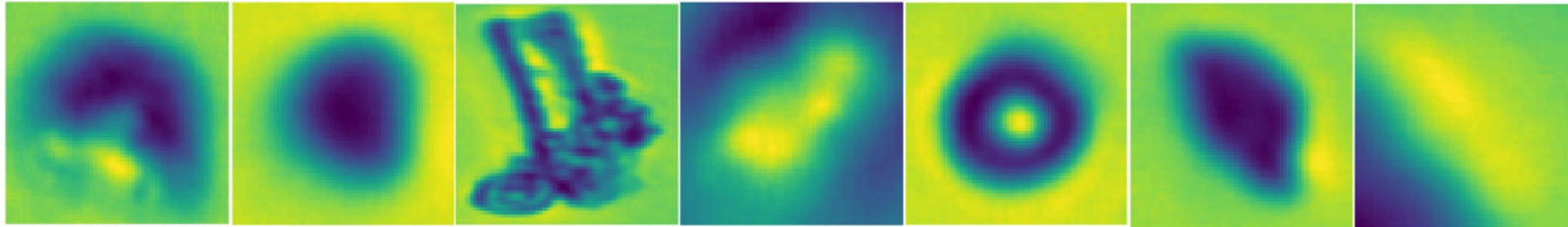


Findings

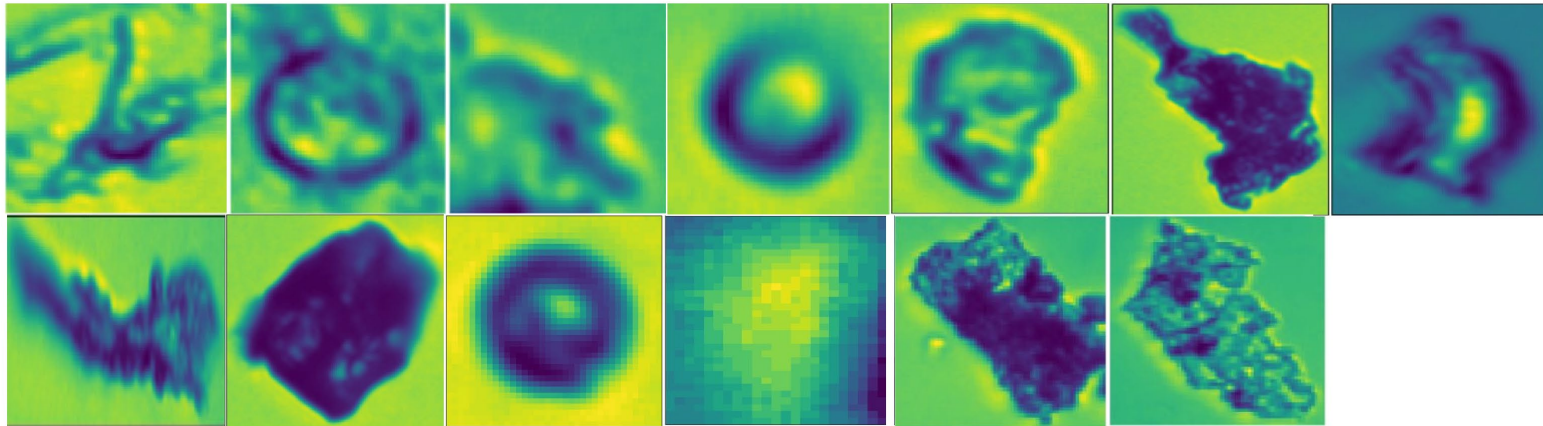
Our results

20 Images

Main categories

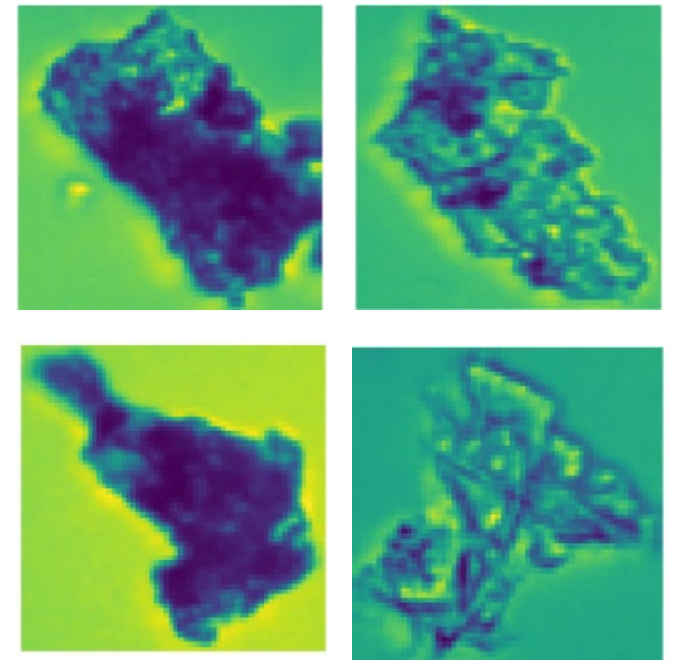
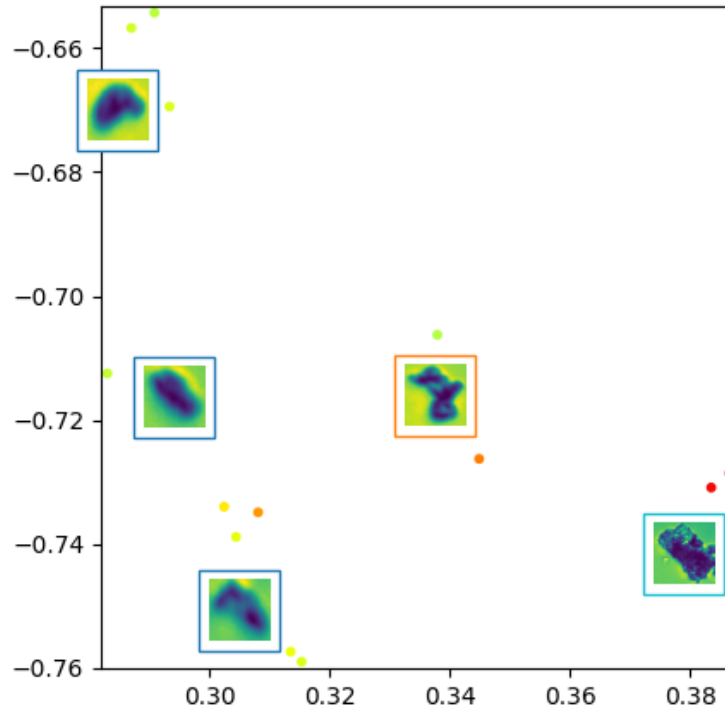
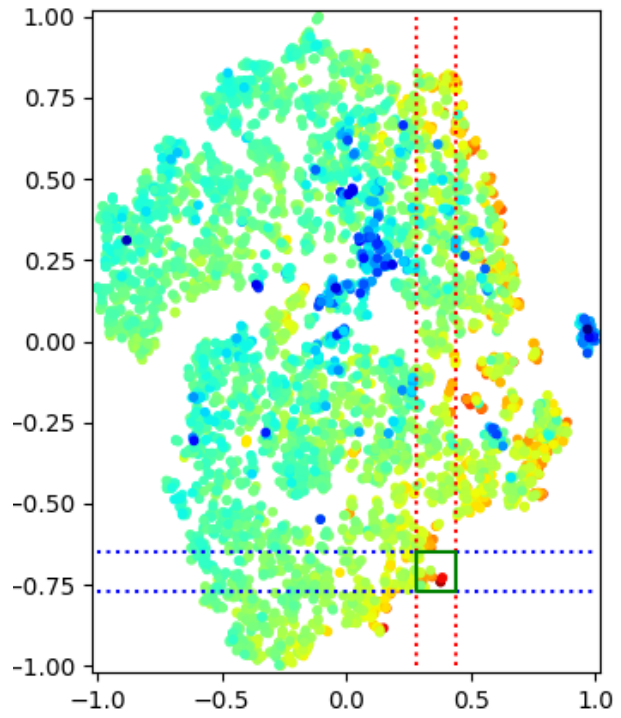


Outliers



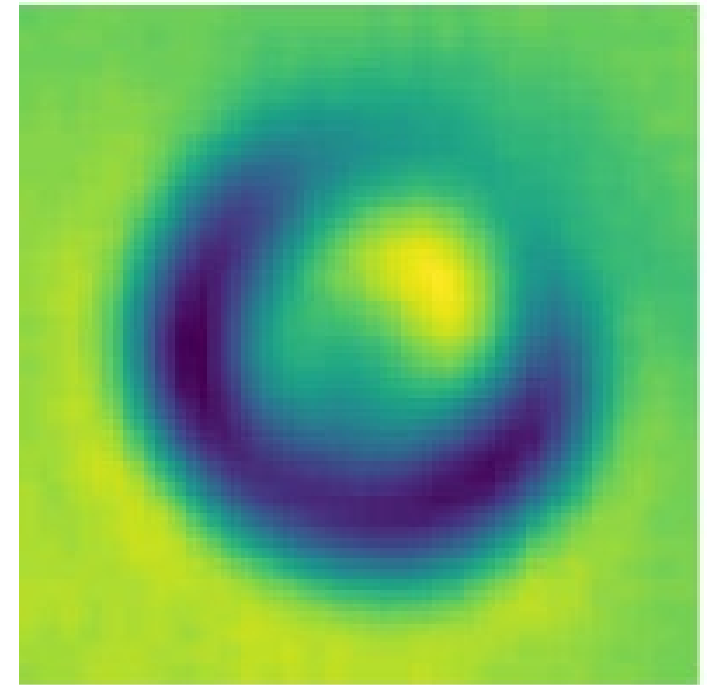
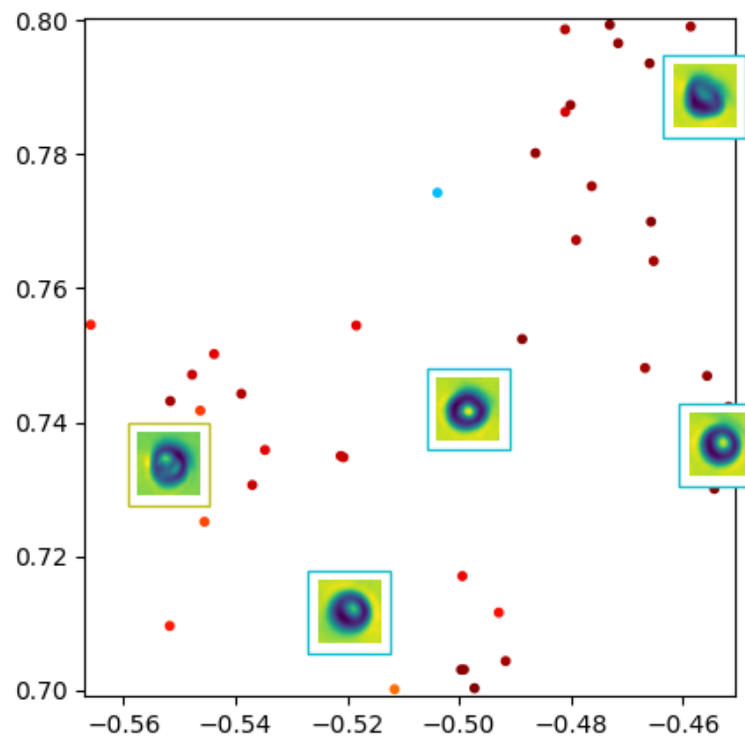
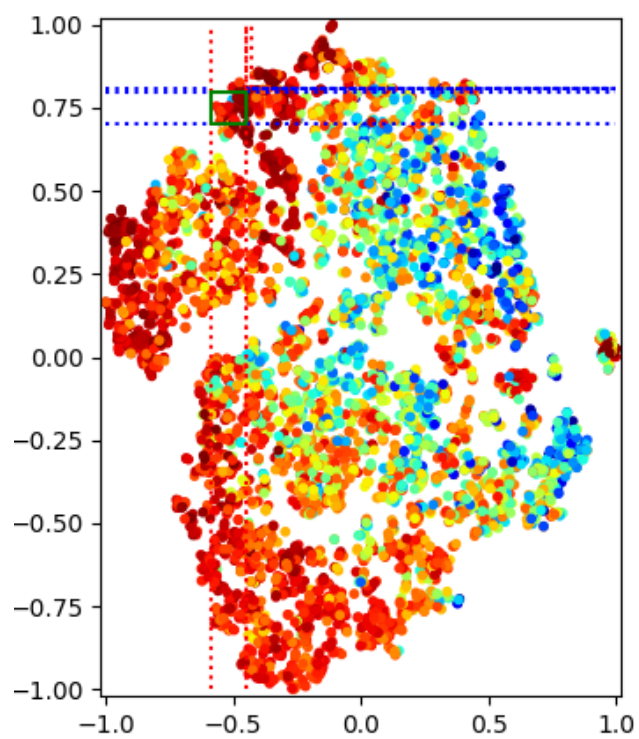
Special 1

Separate area of high-resolution images similar area #5 (Outlier)



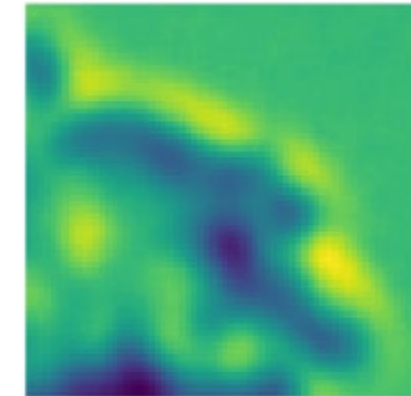
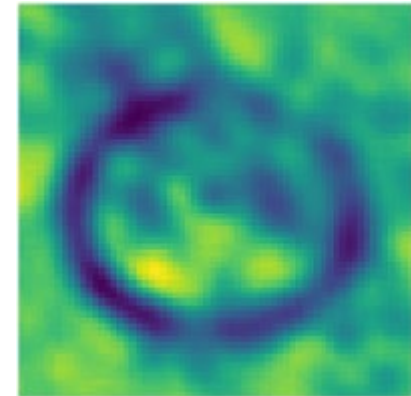
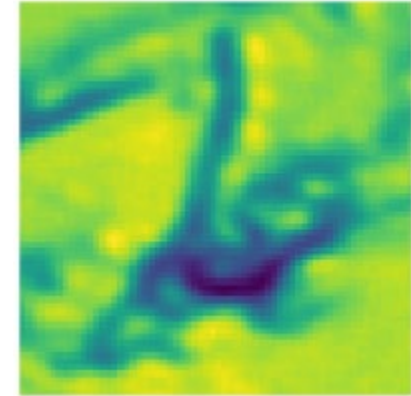
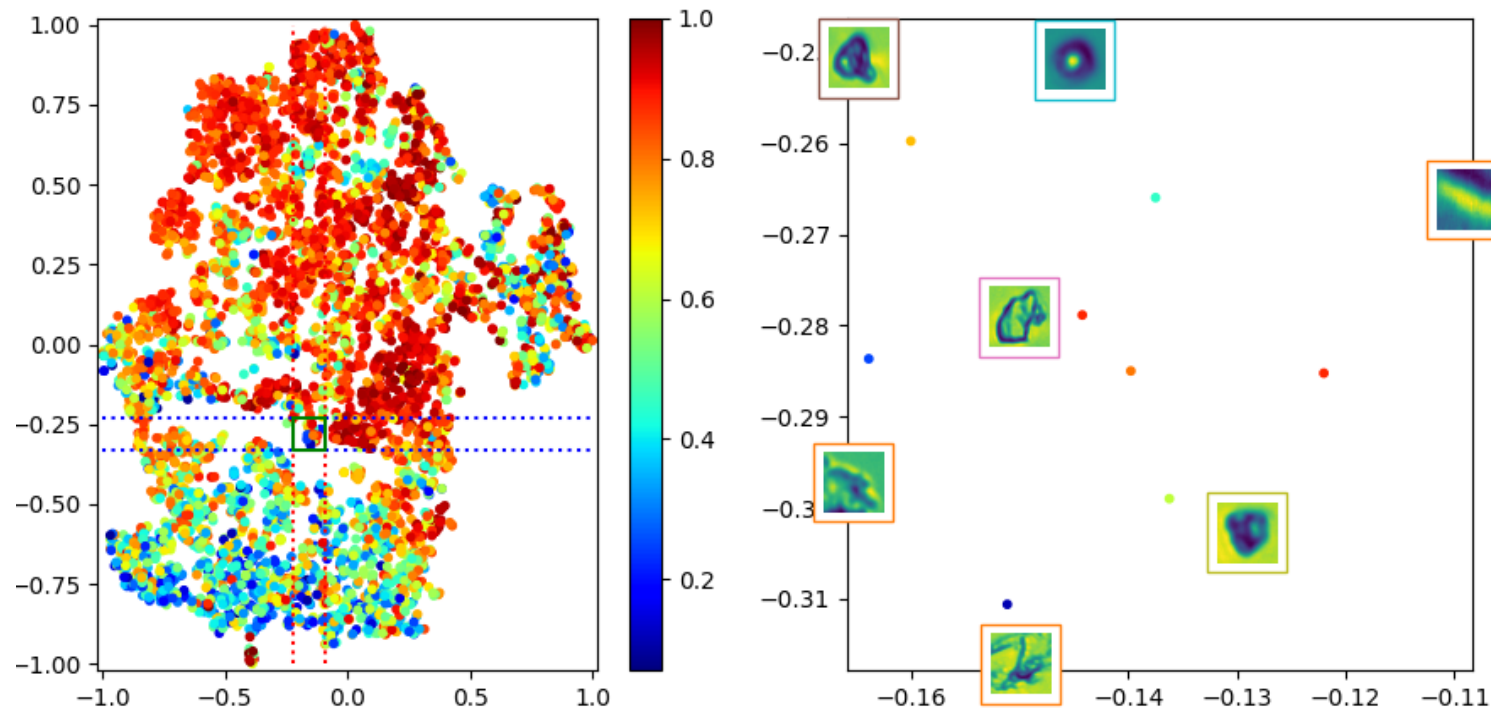
Special 2

Low circularity image surrounded by ones with high circularity. (Outlier)



Special 3

Low circularity images surrounded by ones with high circularity. (Outlier)
Unusually deform

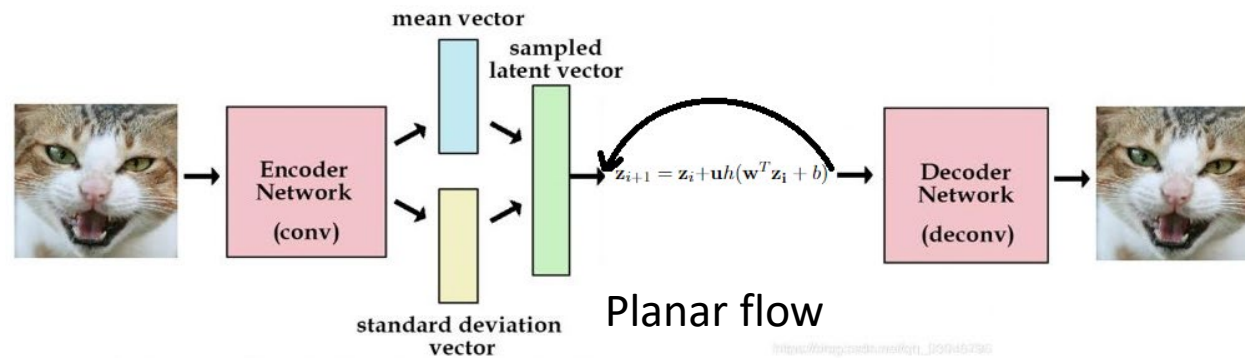
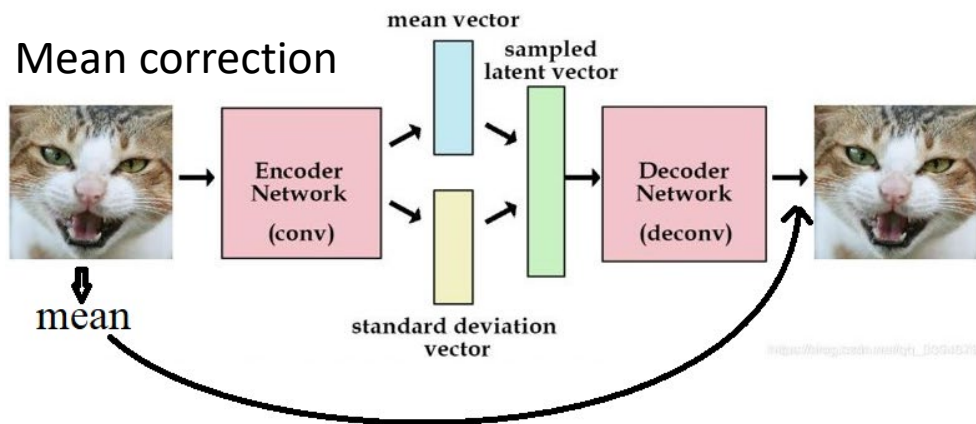


Further work

- Try implementing auto encoders with added scalar output.
- Use permutation importance on scalar parameters in classifier to reduce the number of parameters.
- Work further with the dataset to see if more balanced data for training can be created.

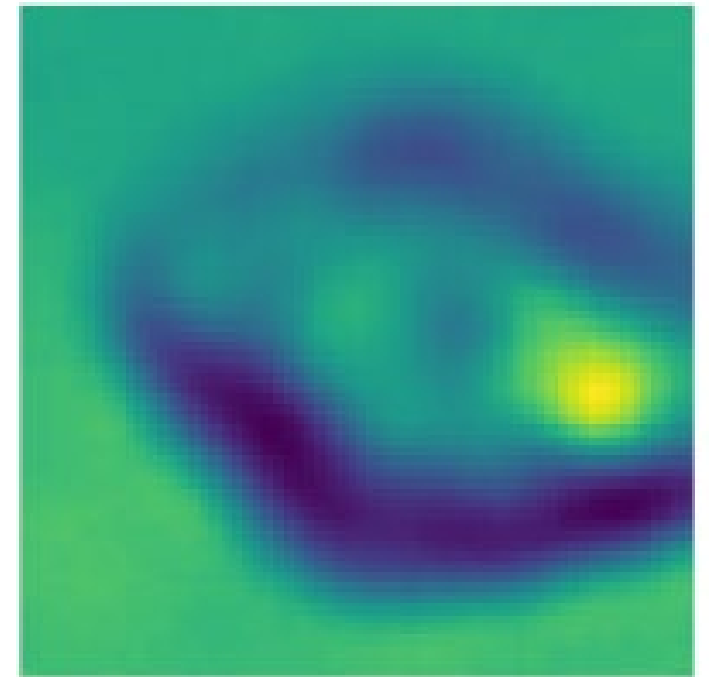
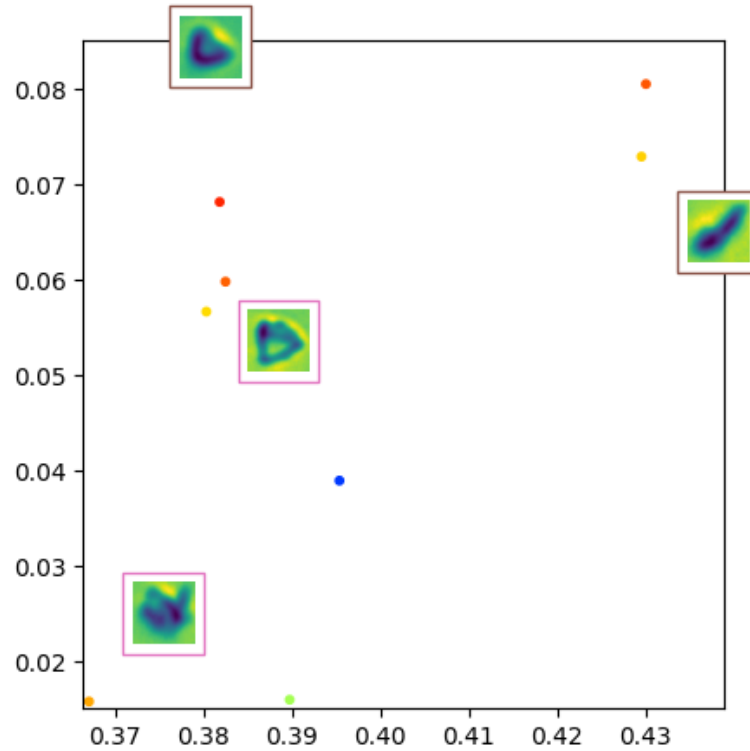
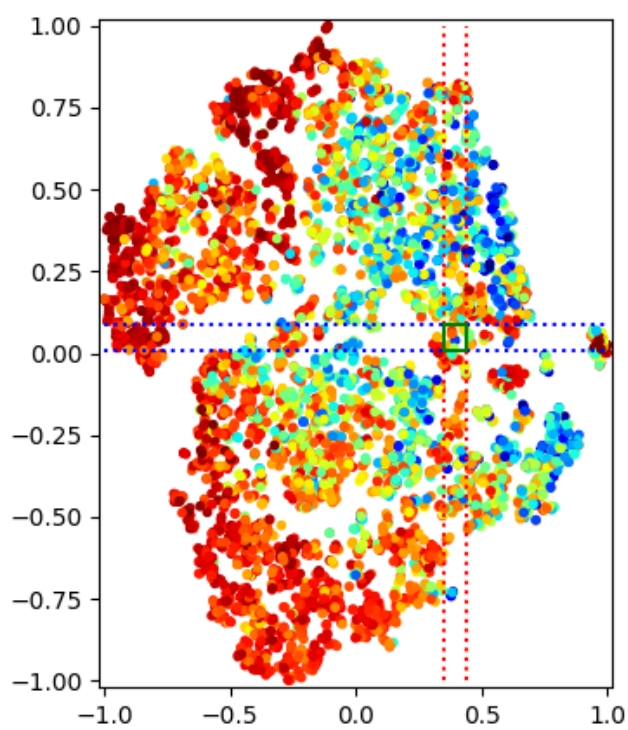
Conclusion

- Using MNIST we were able to reconstruct all classes on partially trained classifiers and from latent extracted from AutoEncoders and VAEs
- Using resnet18 we were able to build models that can make representations of Peruvian ice core images.
- Encountered and overcame problems while training VAEs.
- Fitting less to the artificial set seemed to help the VAE classes.
- Tried clustering and different (V)AE architectures.
- Were able to use representations and dimensionality reduction to identify main classes and outlier images.



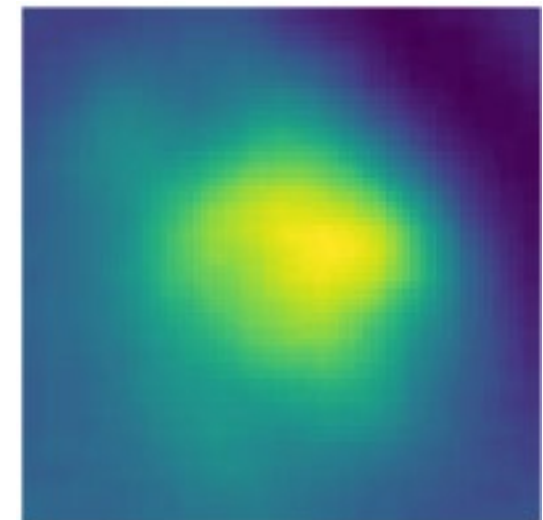
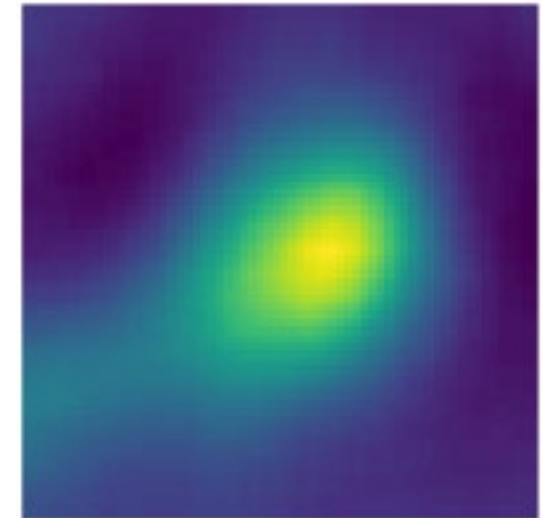
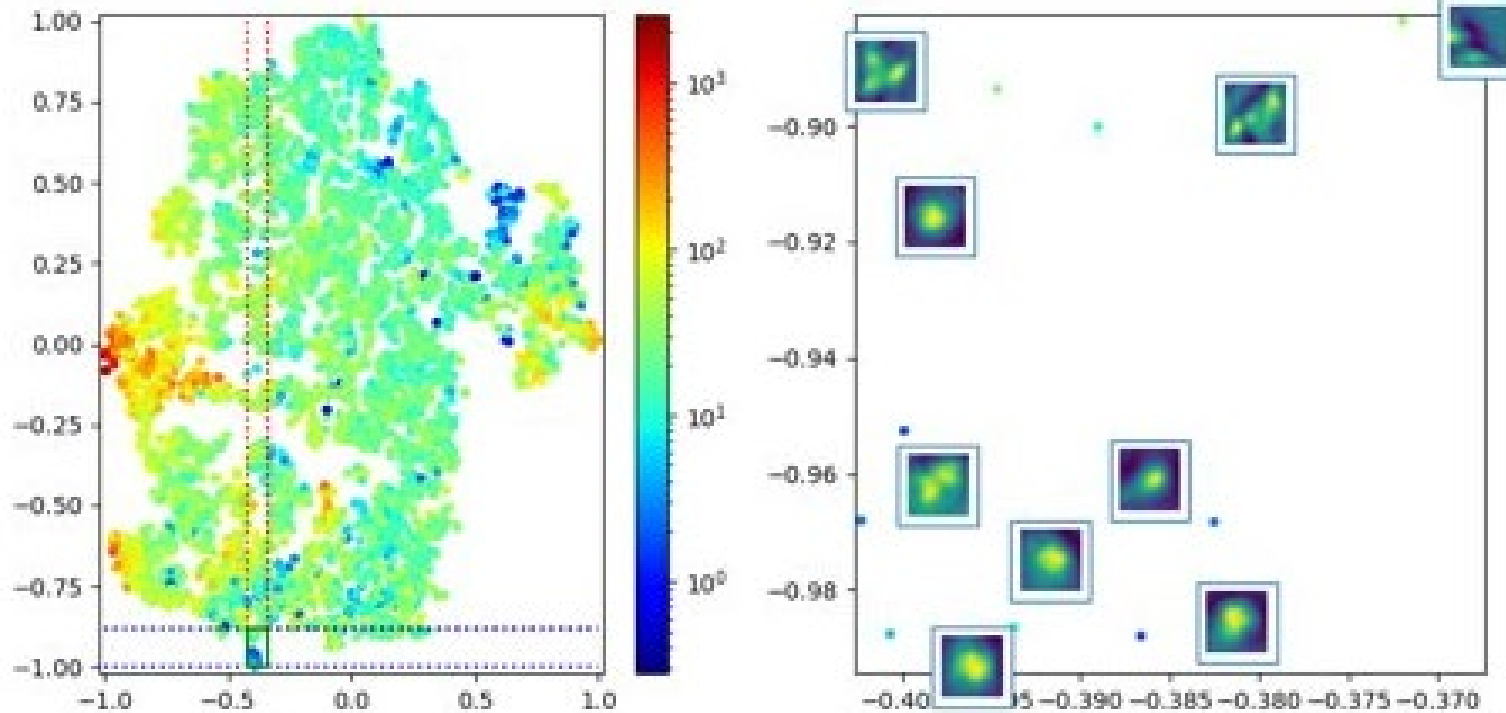
Special 4

Low circularity image surrounded by ones with high circularity. (Outlier)



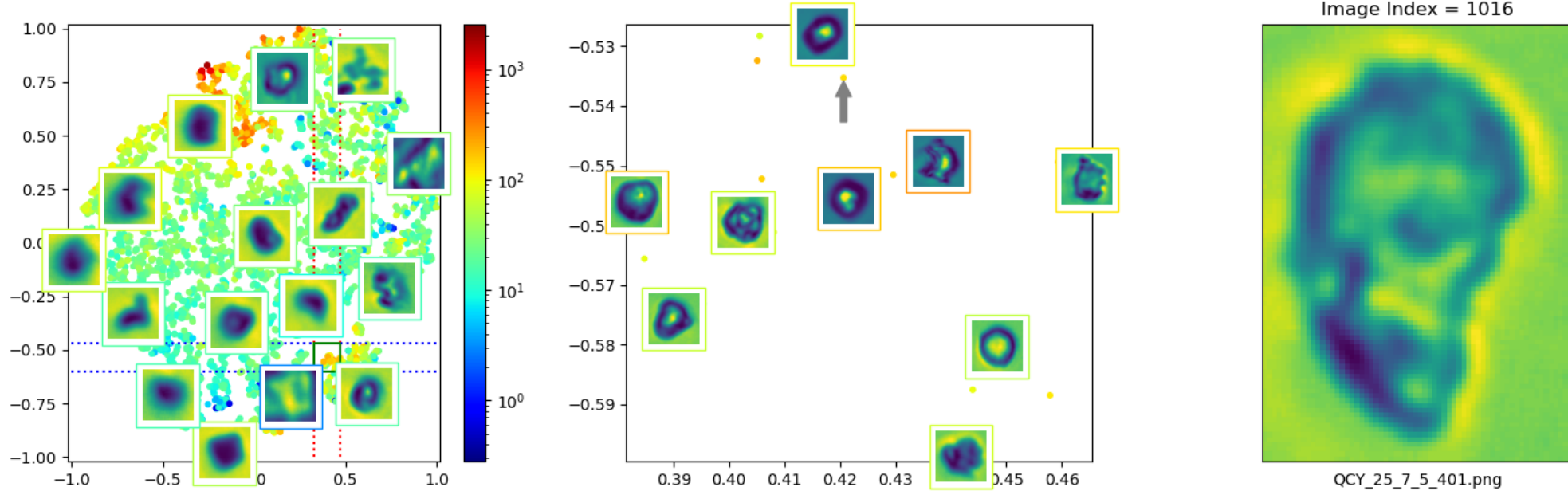
Special 5

Local area of small image "light blobs"



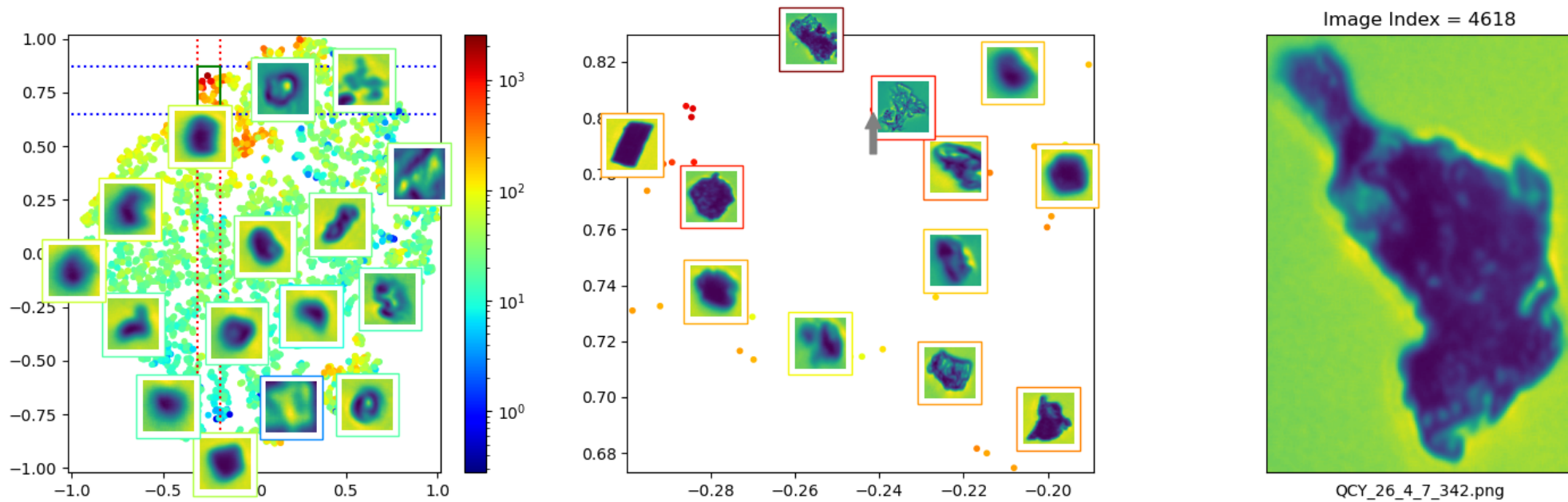
QCY_25_7_5_401

Why?



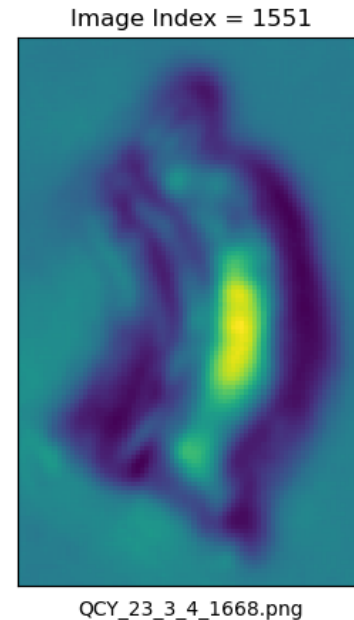
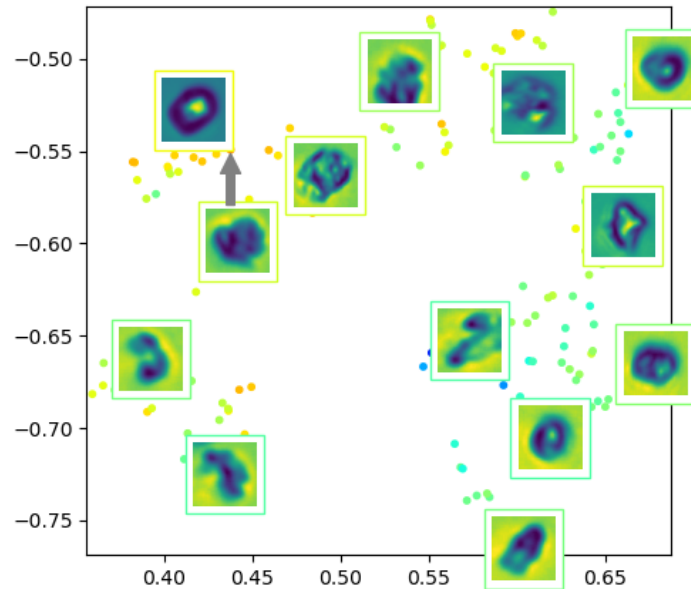
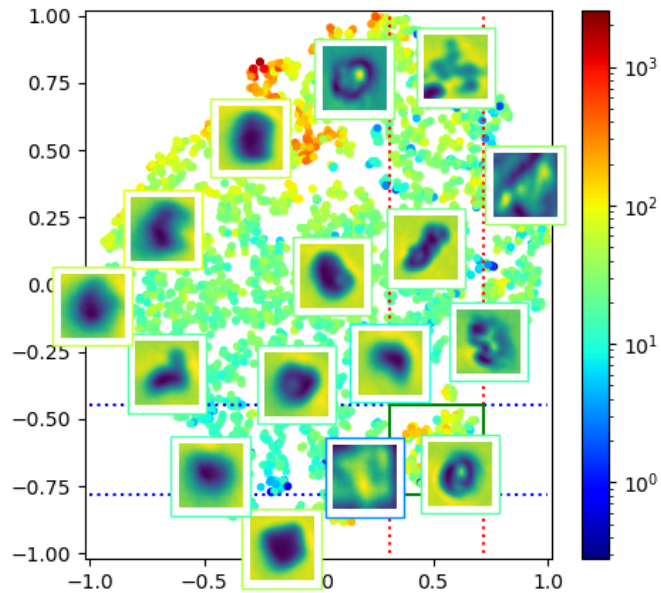
QCY_26_4_7_342

Why?



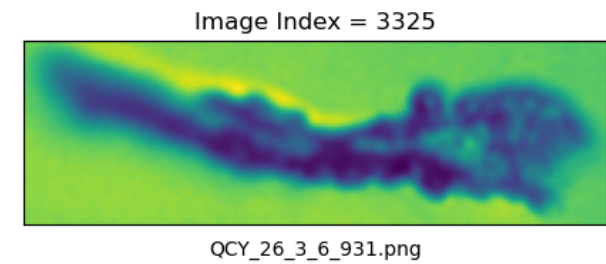
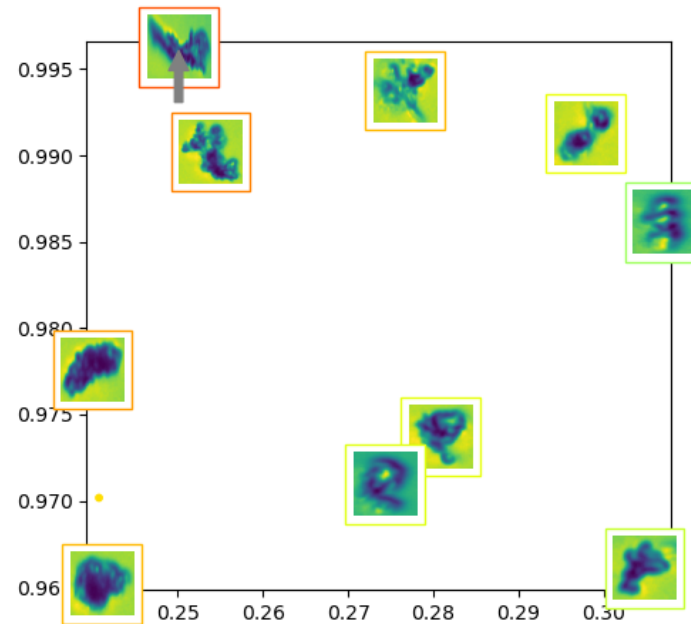
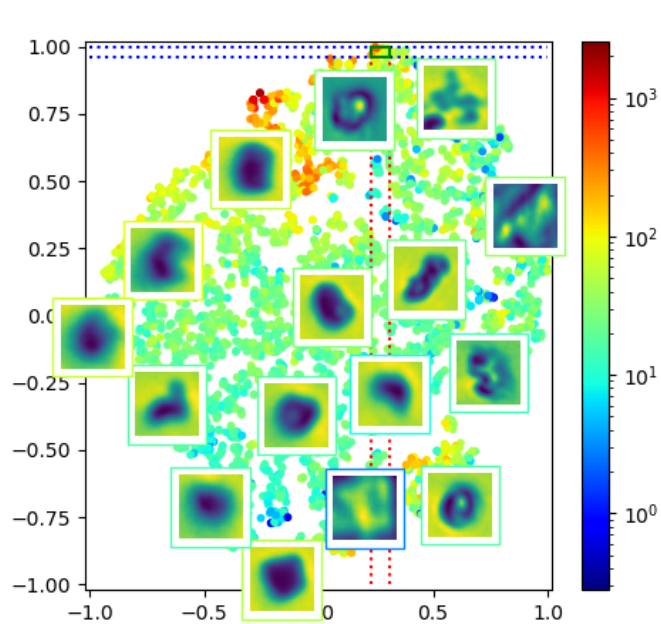
QCY_32_3_4_1668

This looks to be a separate cluster and the selected images is an example of it



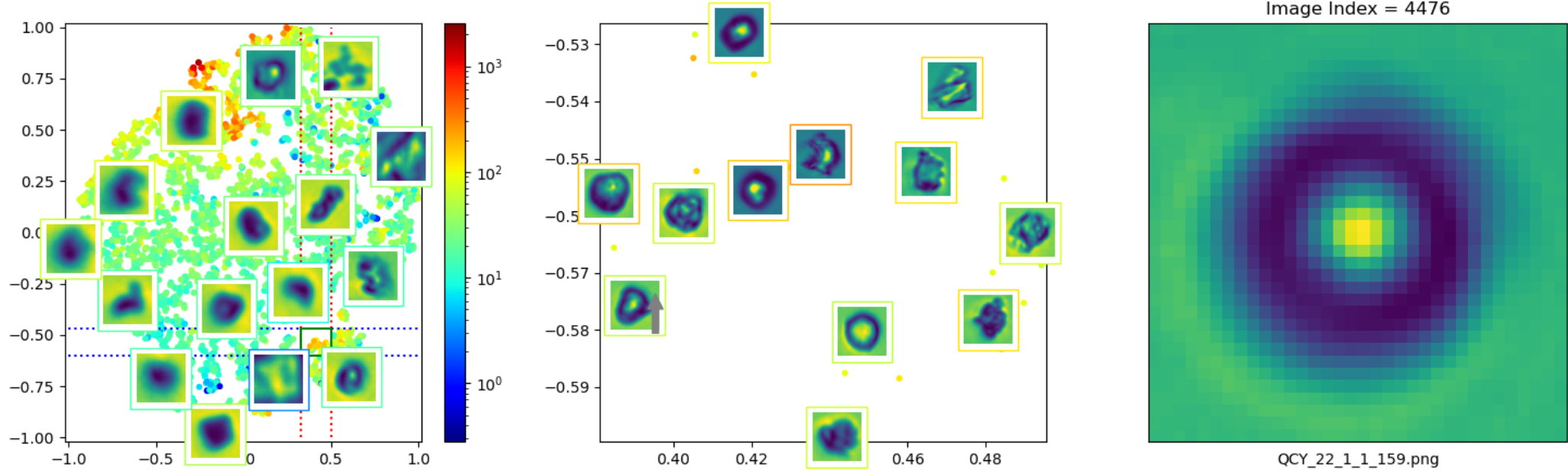
OCY_26_3_6_931

Visually distinct large object. (15th largest object)



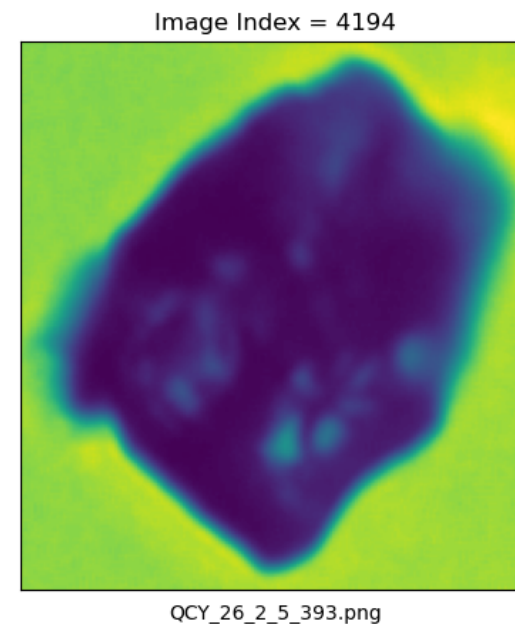
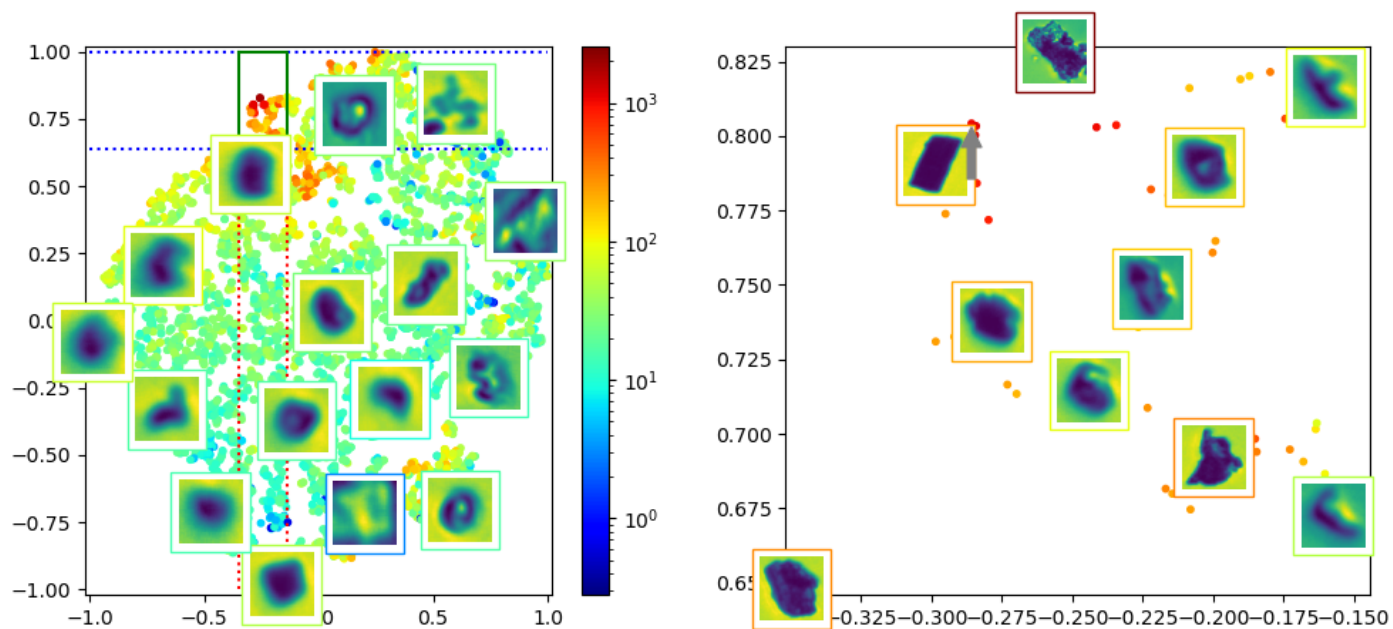
QYC_22_1_1_159

Small image (Outlier)



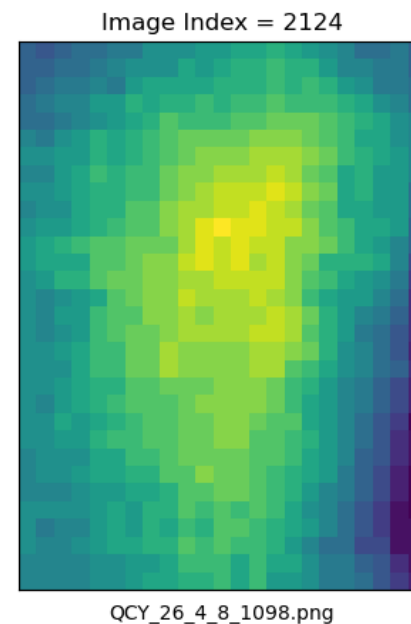
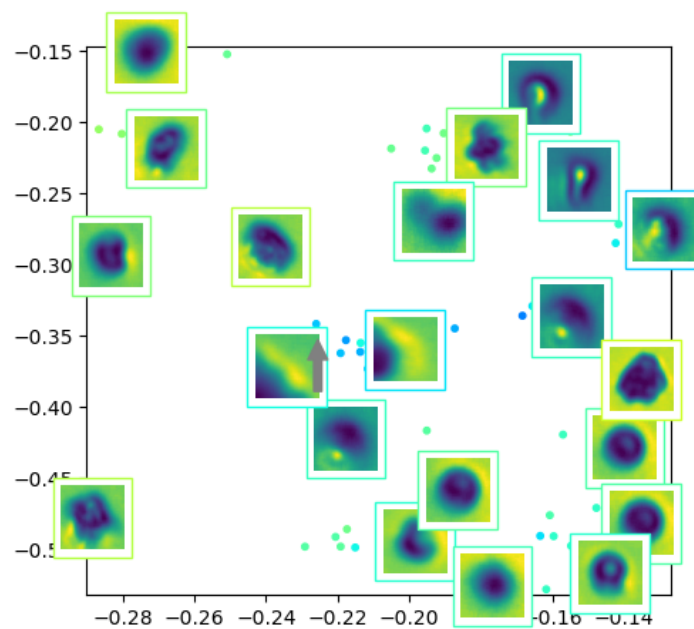
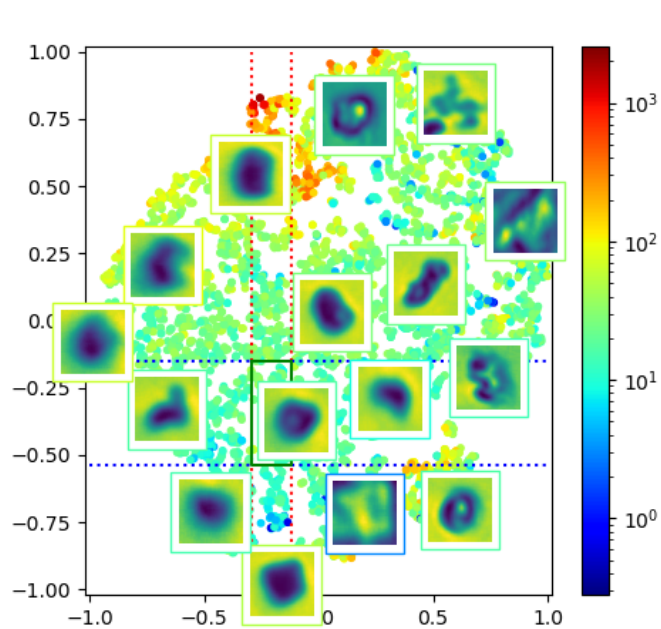
QCY_26_2_5_393

Highest loss for on VAE



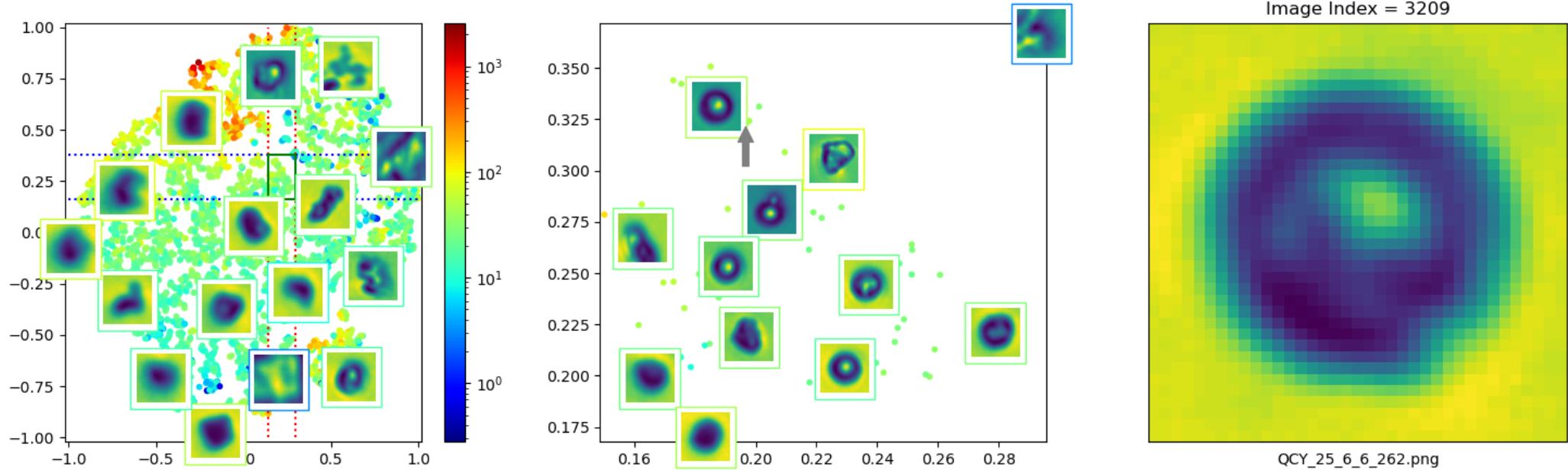
QCY_26_4_8_1098

Second highest loss on VAE



QCY_25_6_6_262

Previously clustered differently

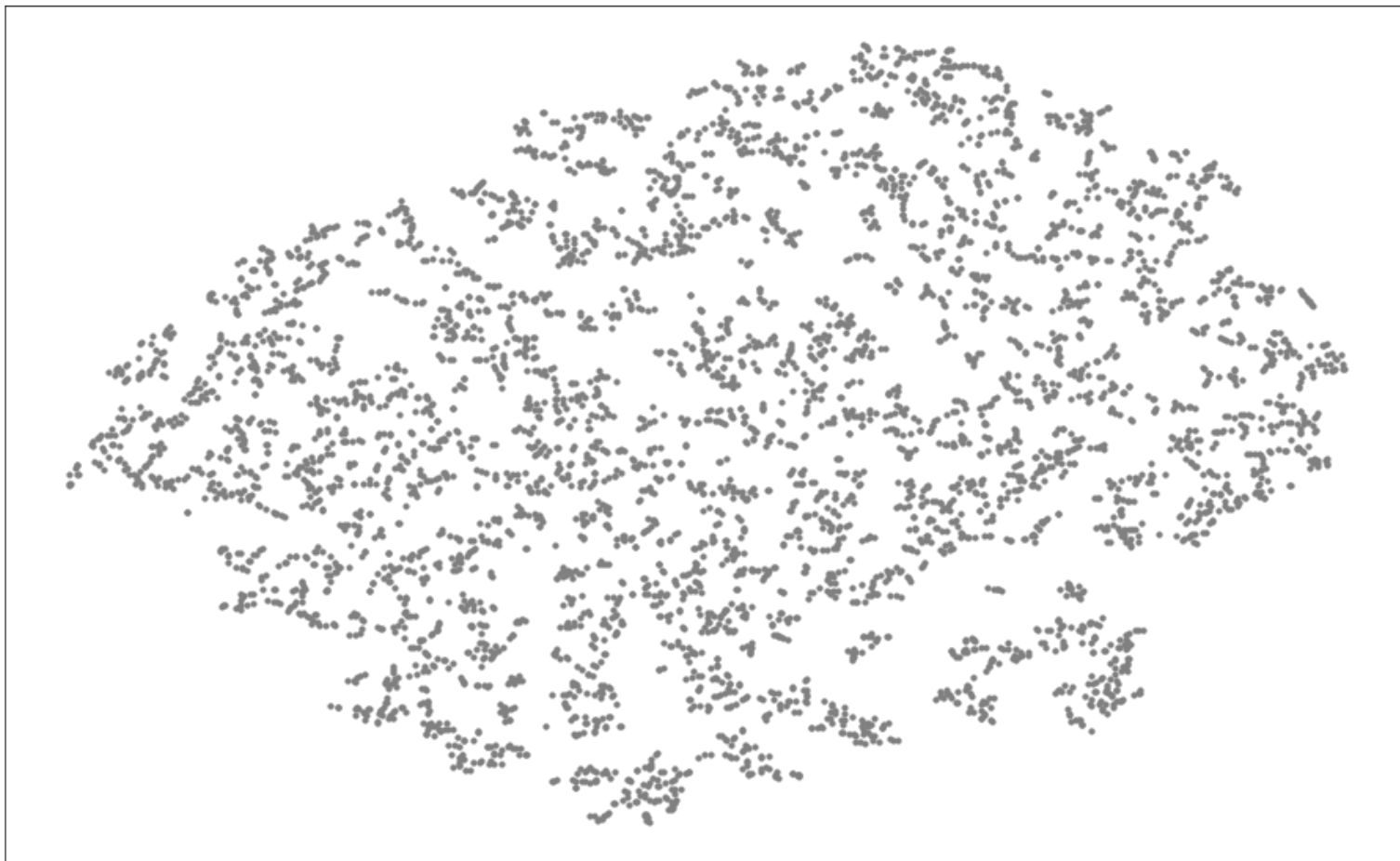


Appendix

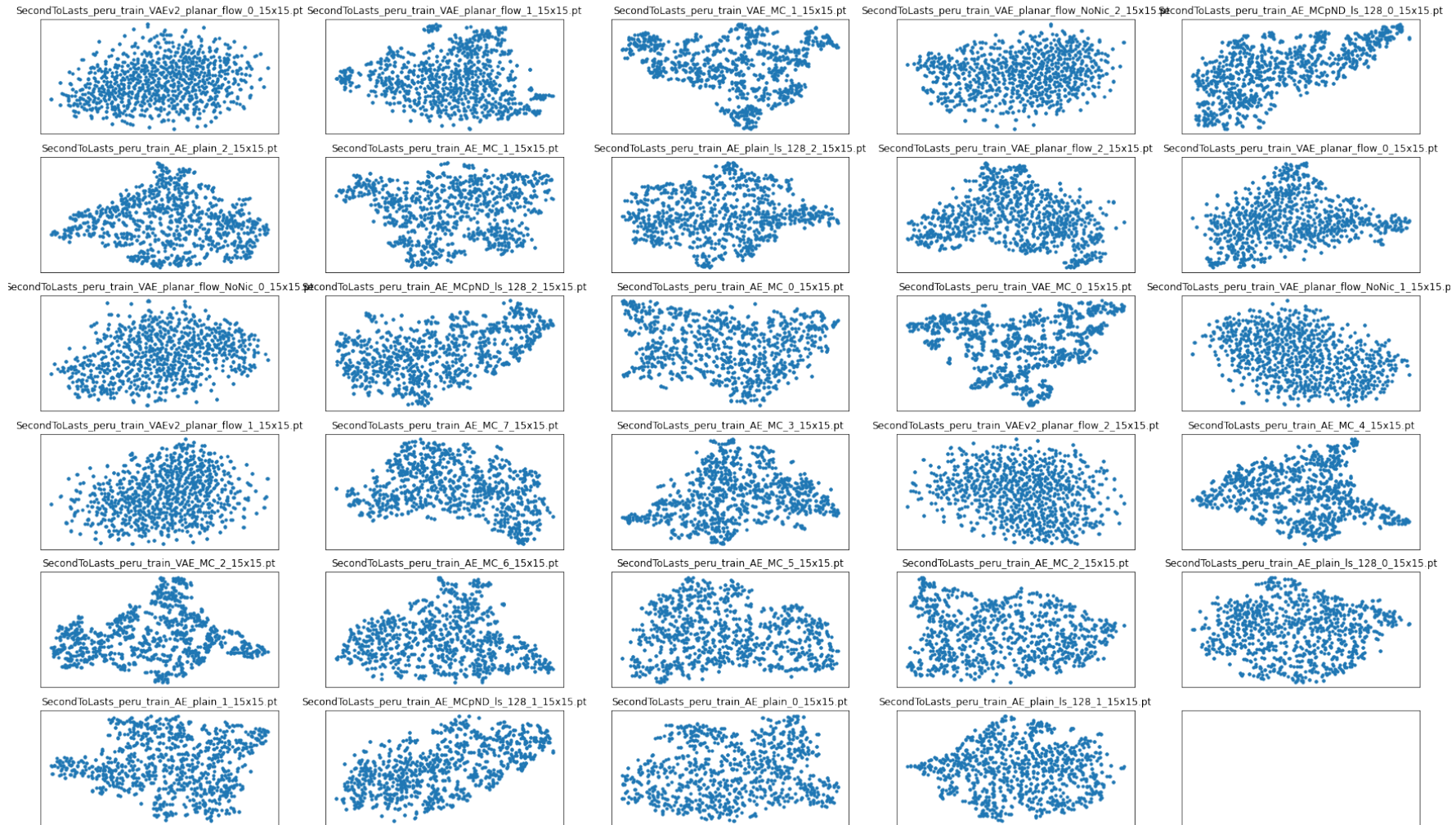
Are you looking for more?

Plain plot

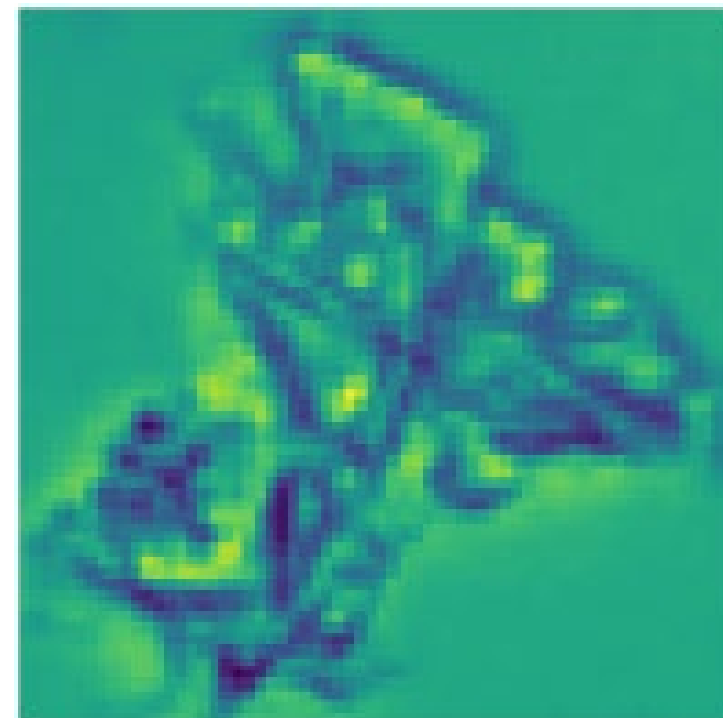
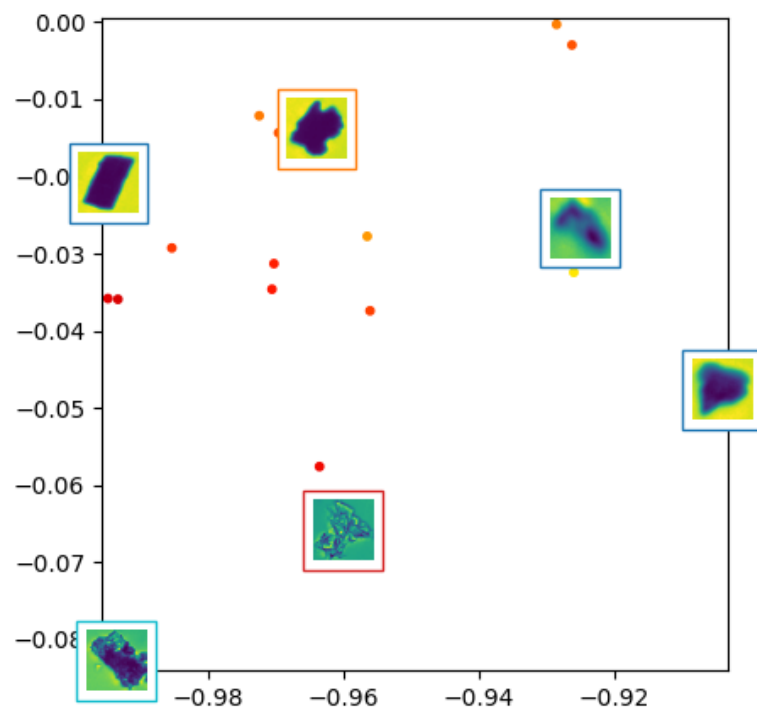
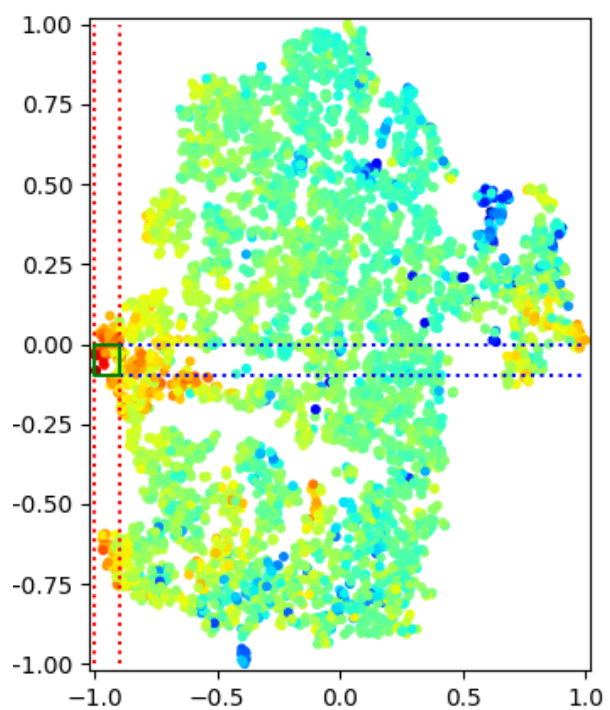
Plot without colors



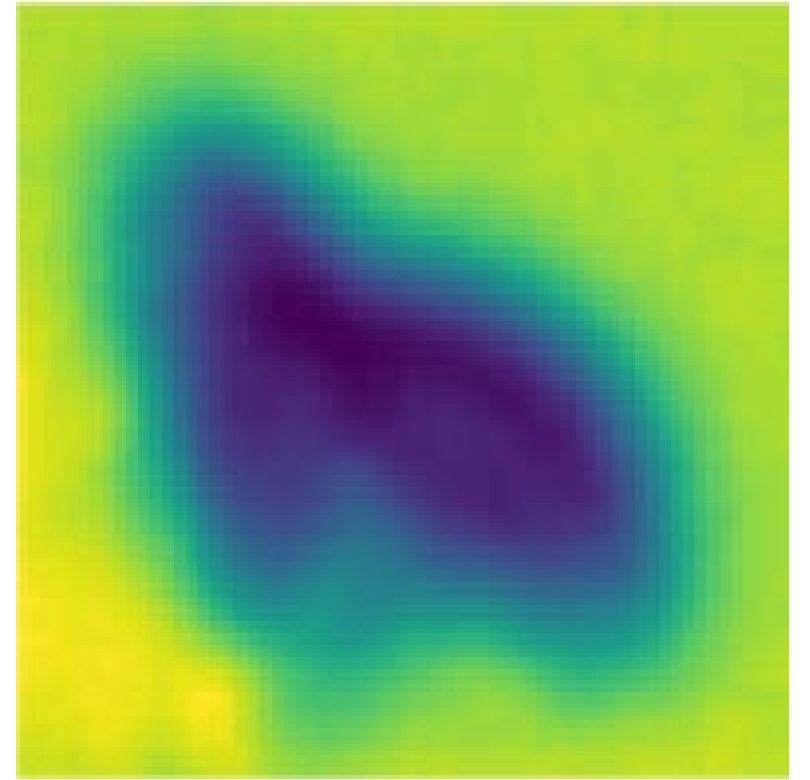
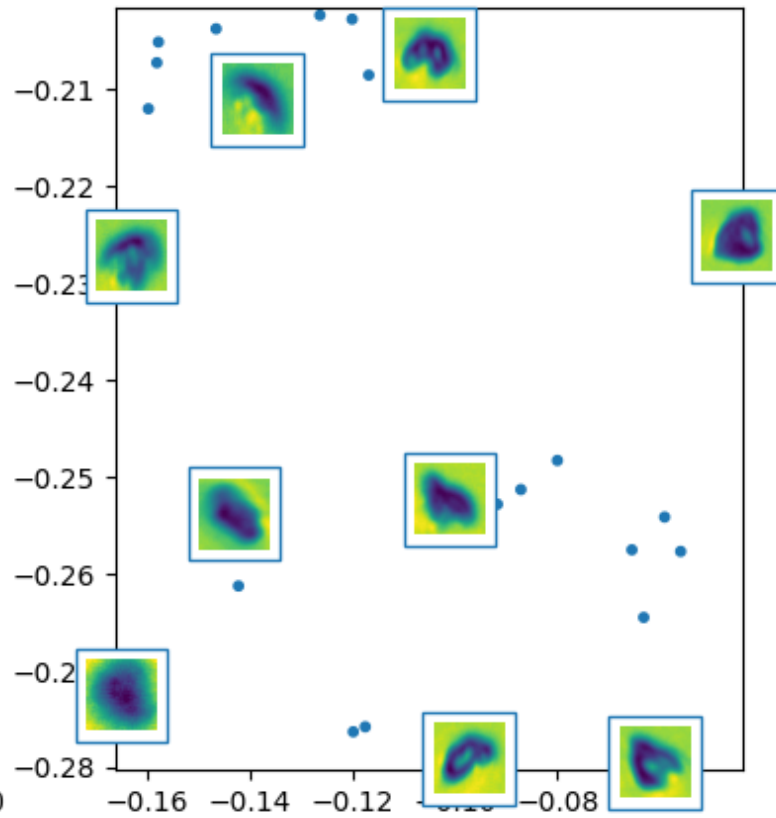
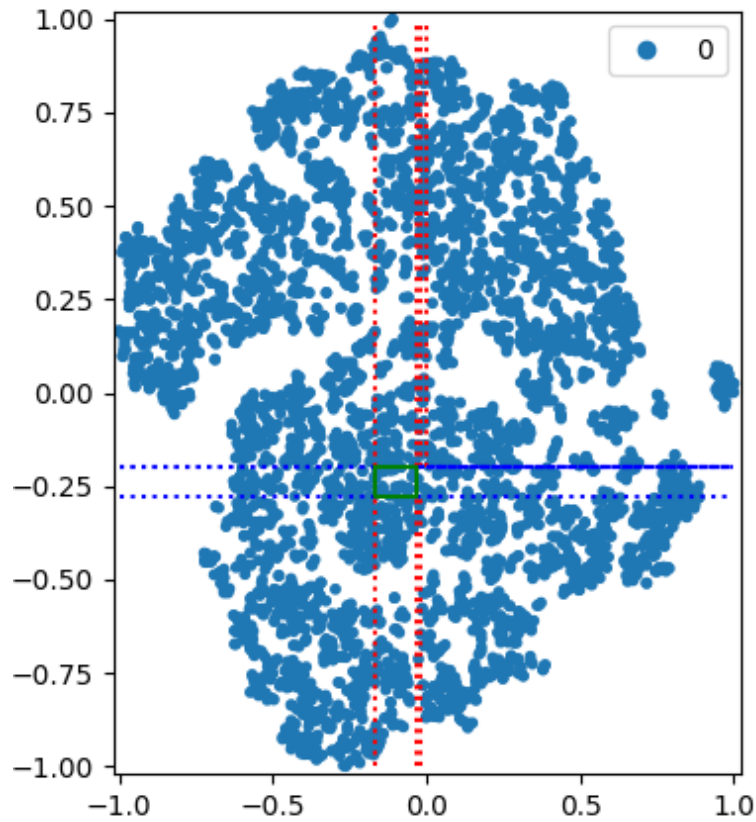
t-Sne for different (V)AEs latent spaces



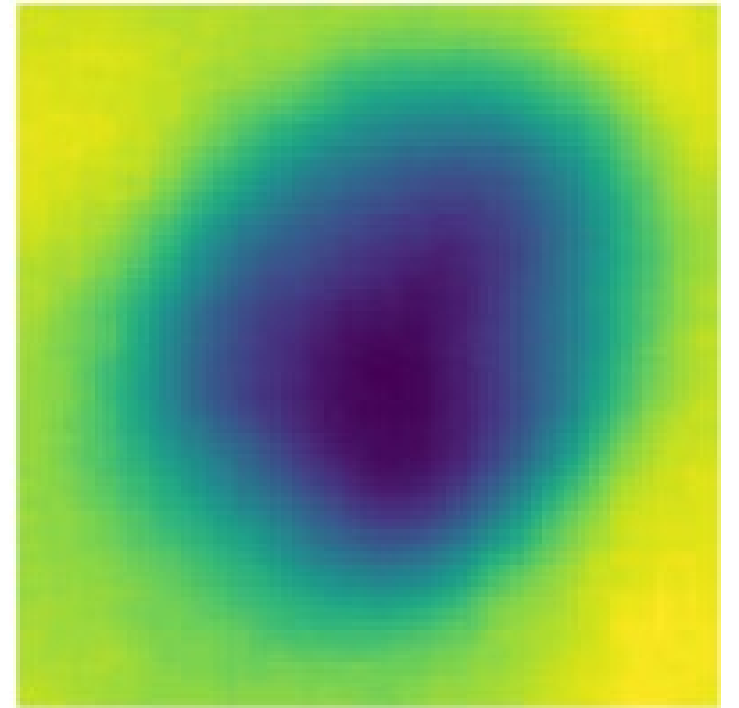
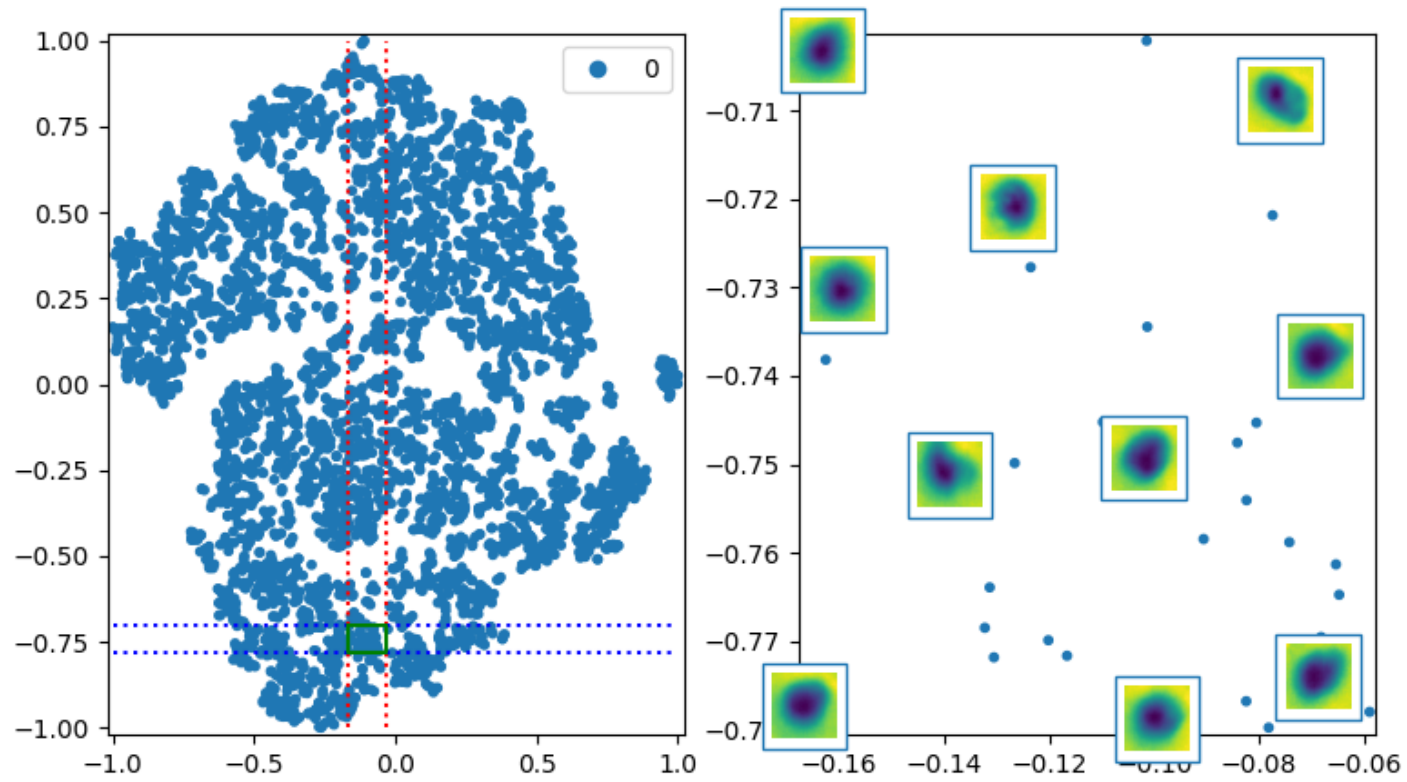
Refound largest images



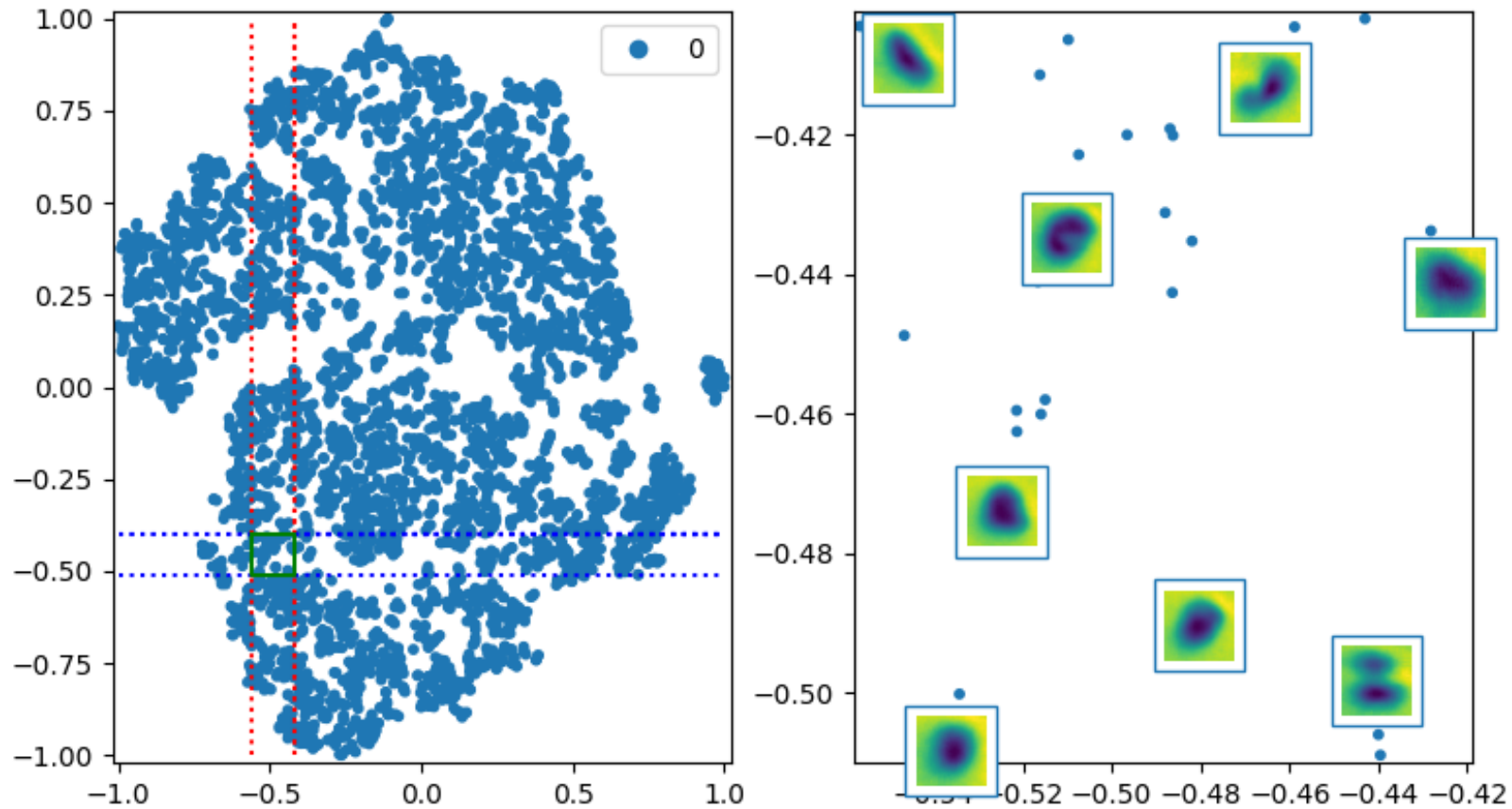
CNN classifier - Area 3



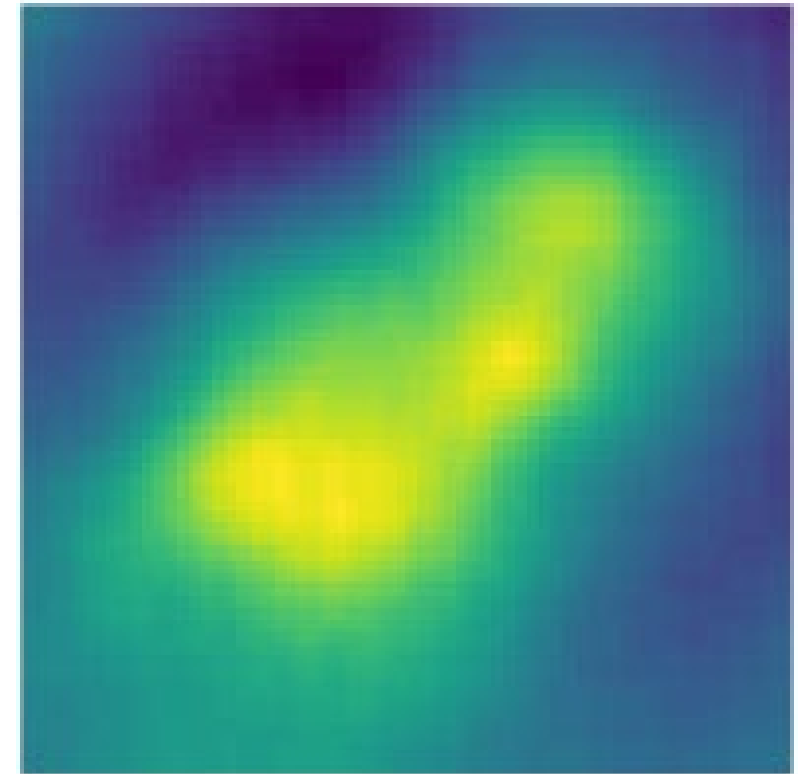
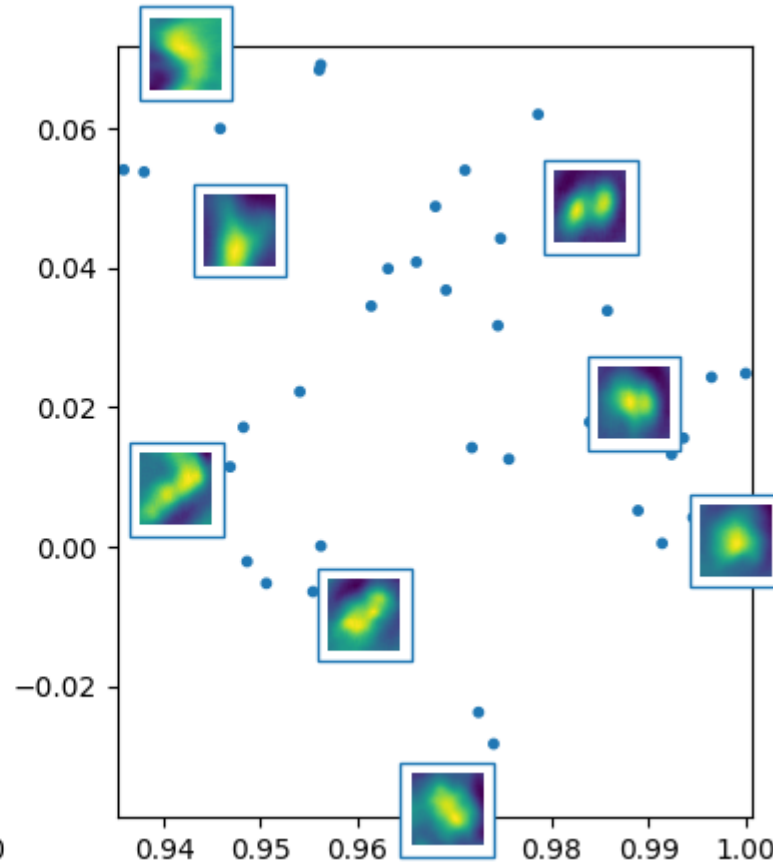
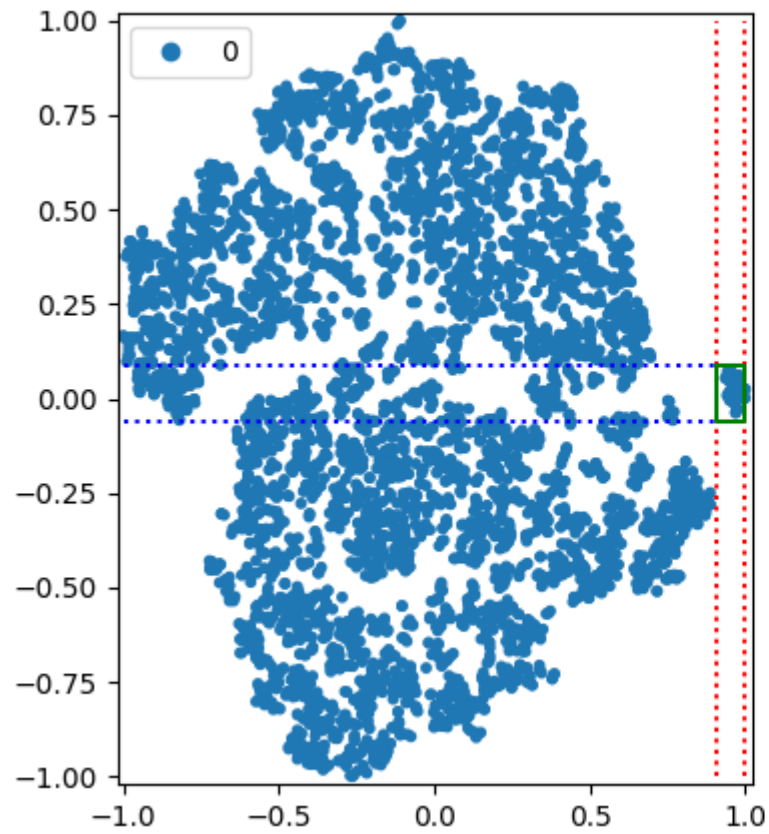
CNN classifier - Area 4



Between area 3 and 4

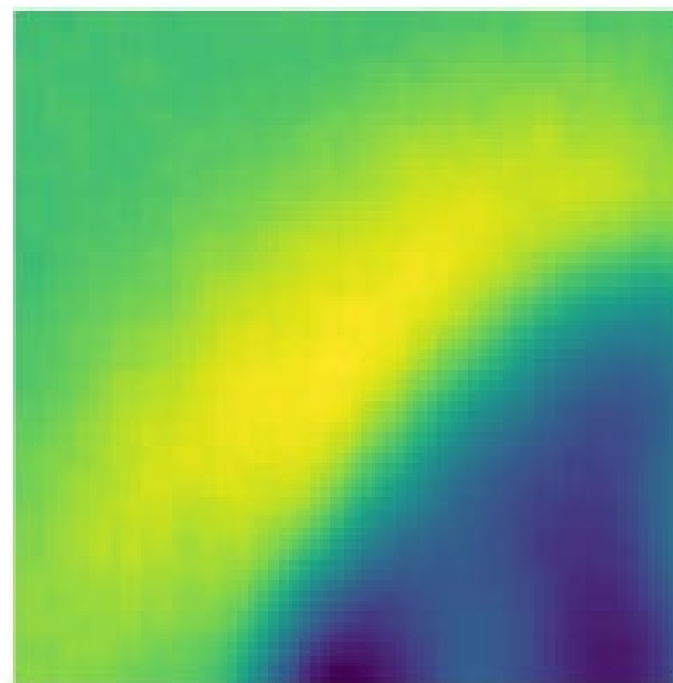
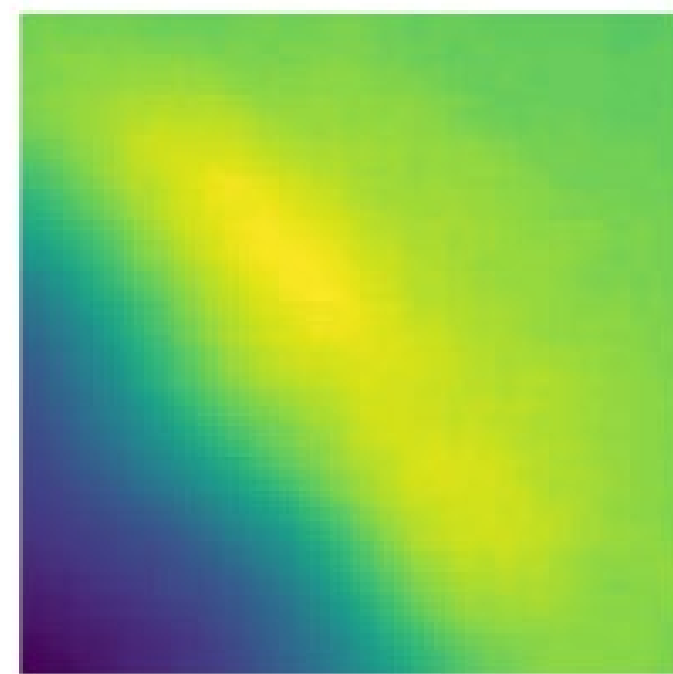
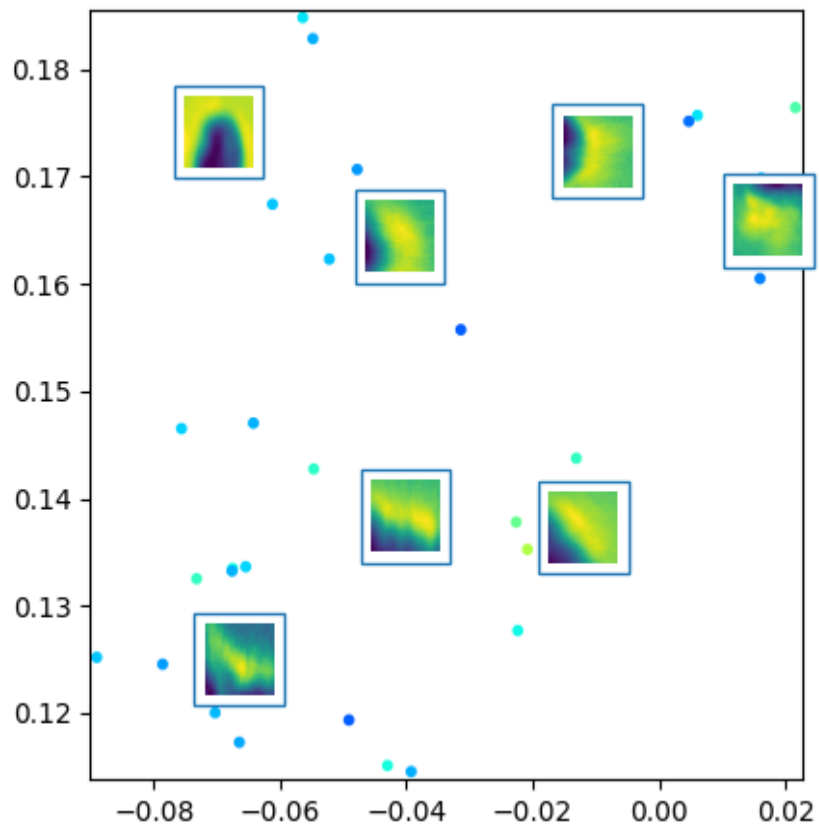
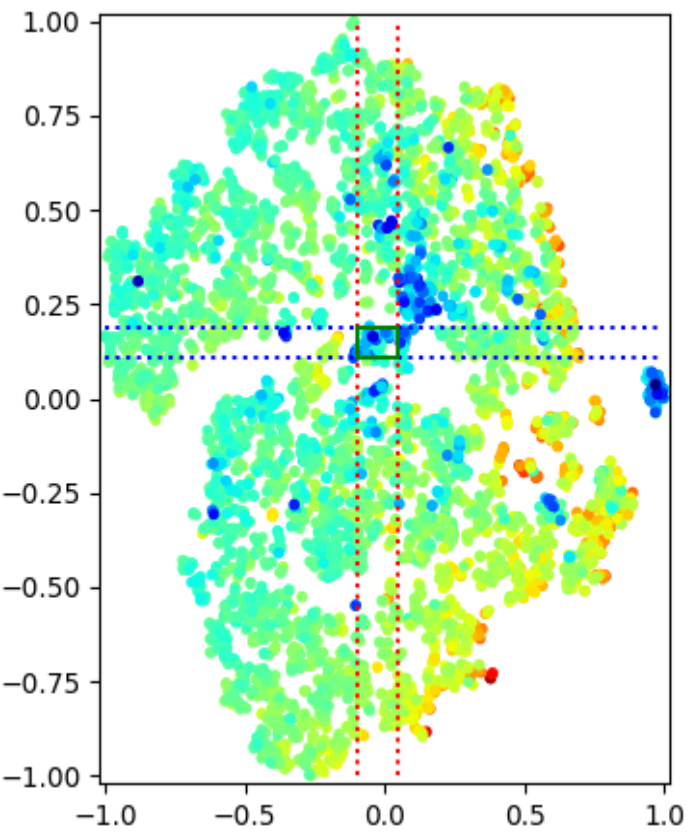


CNN classifier - Area 6

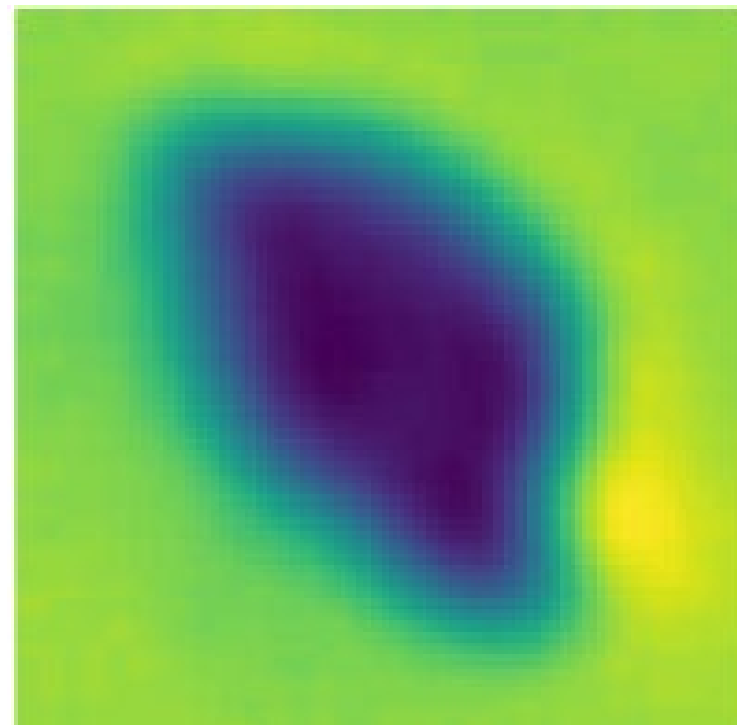
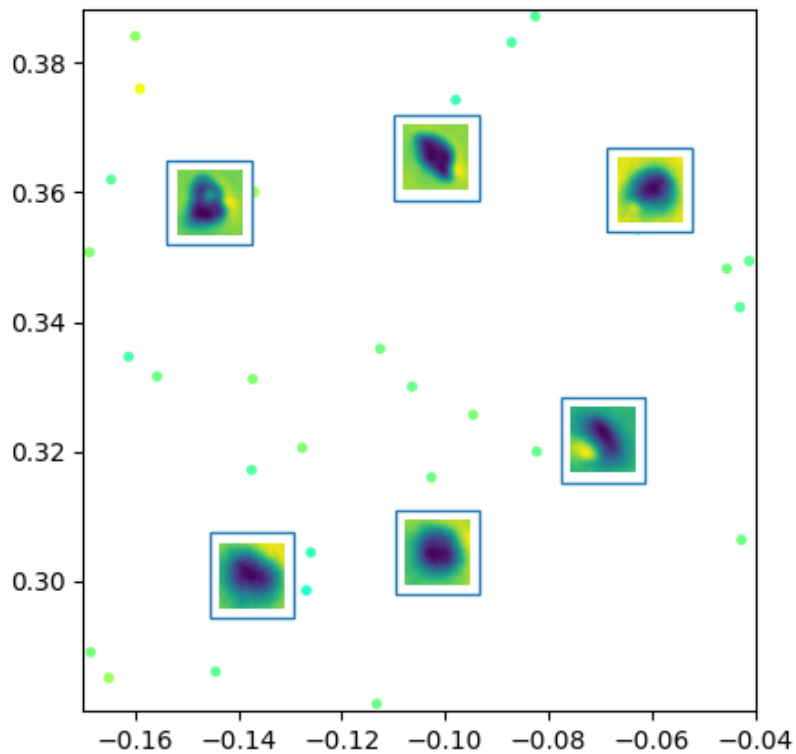
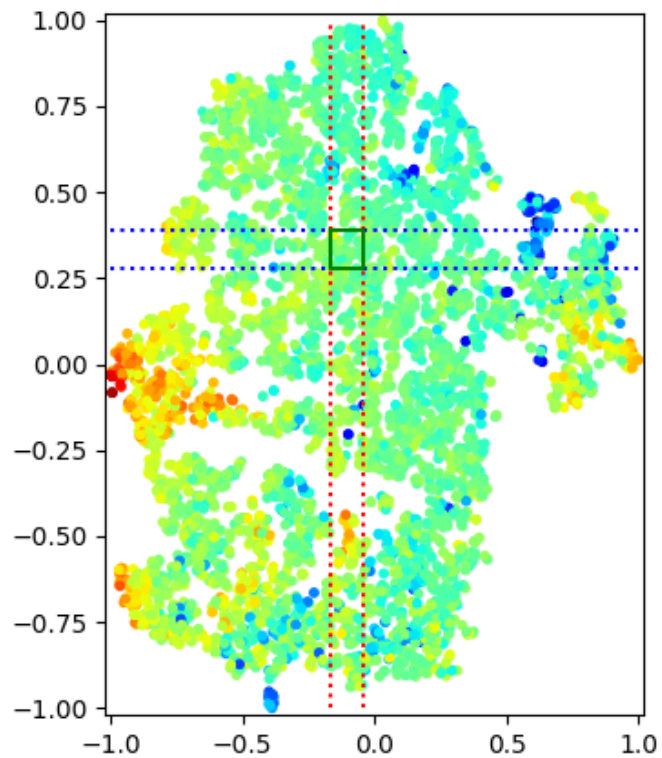


Area 7

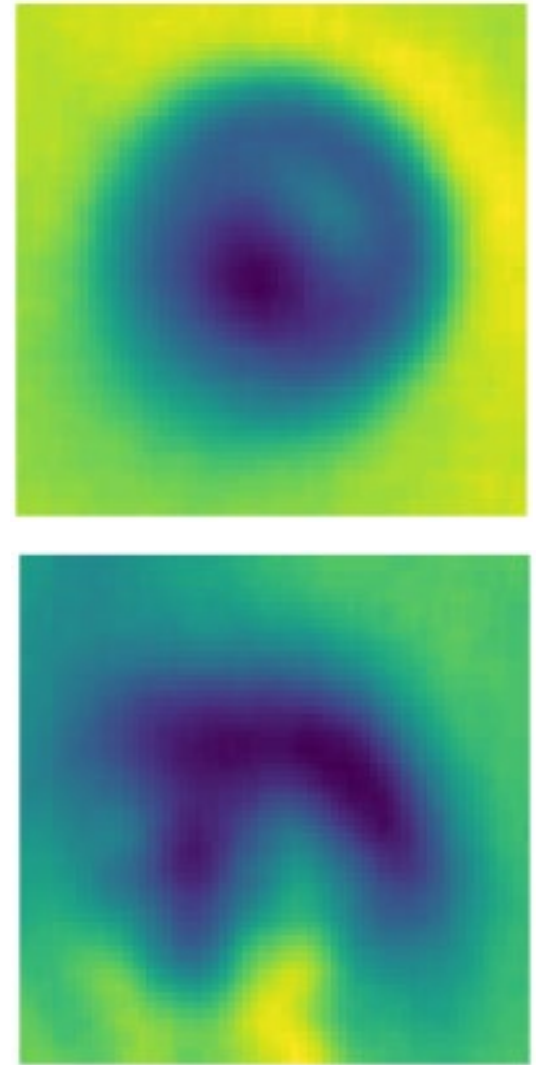
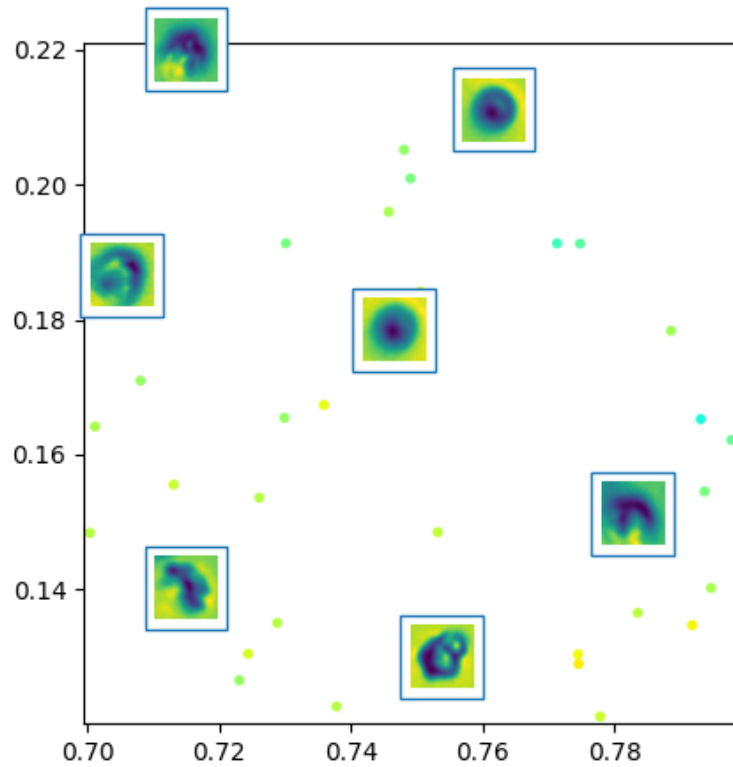
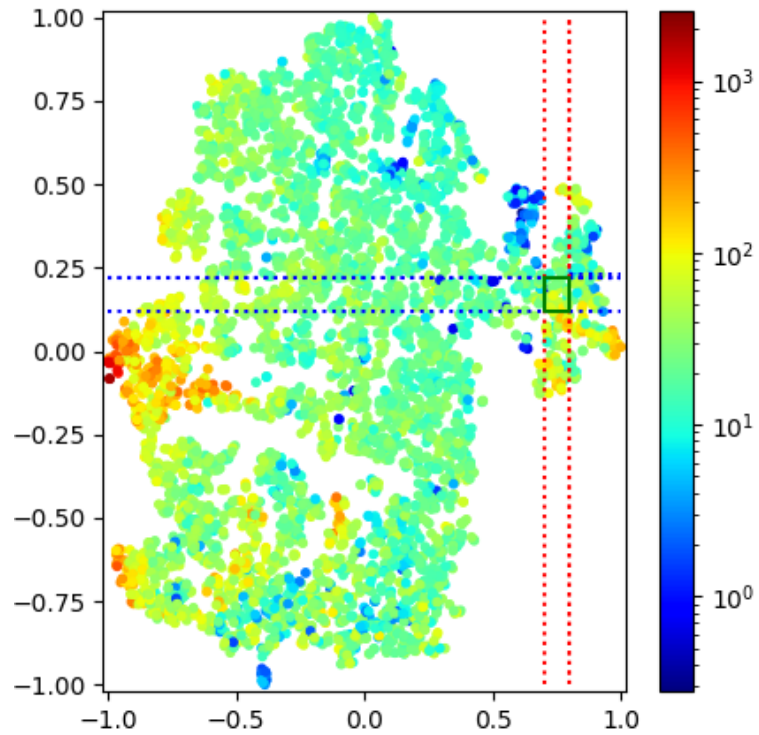
Low area



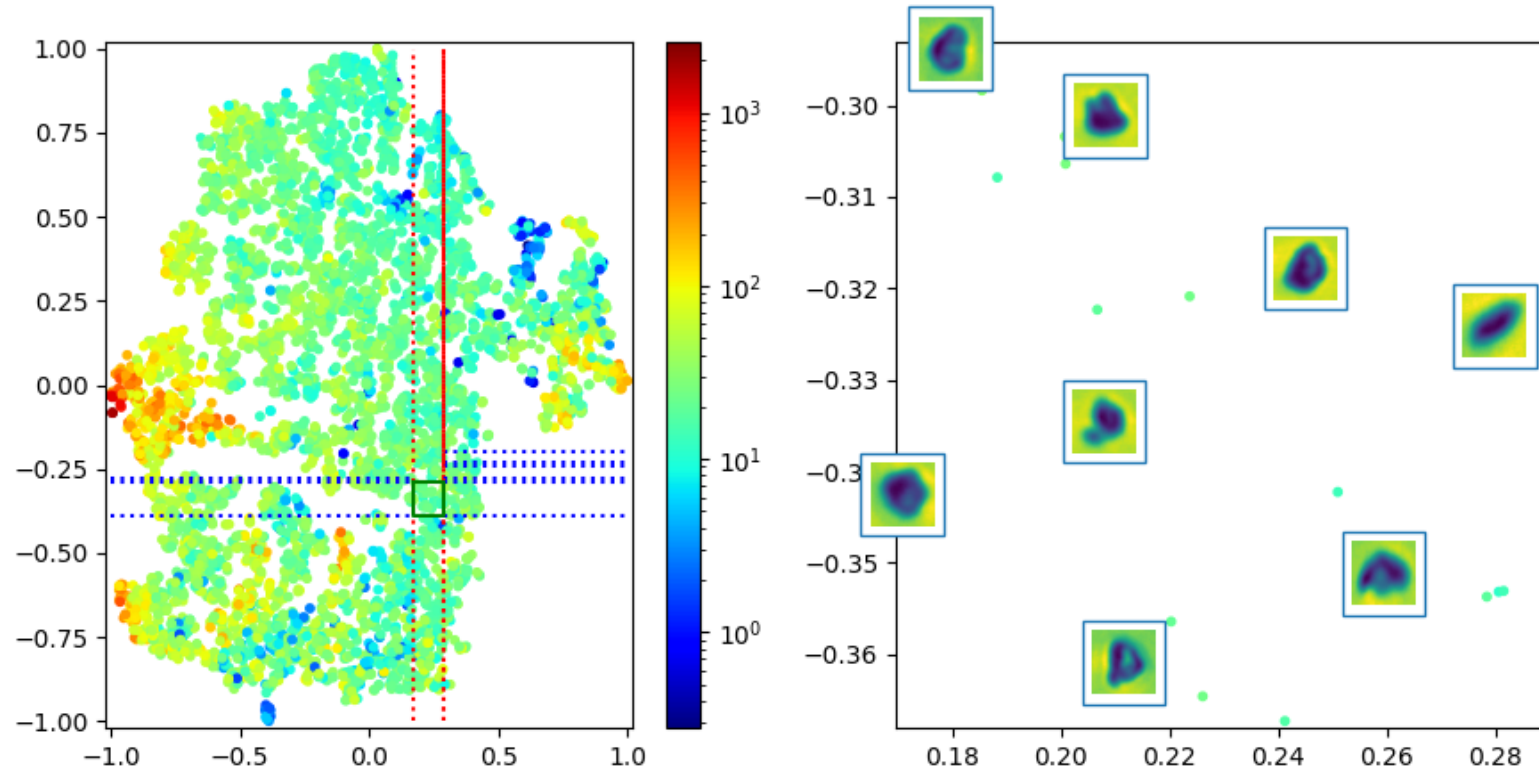
VAE, planar flow Area 1



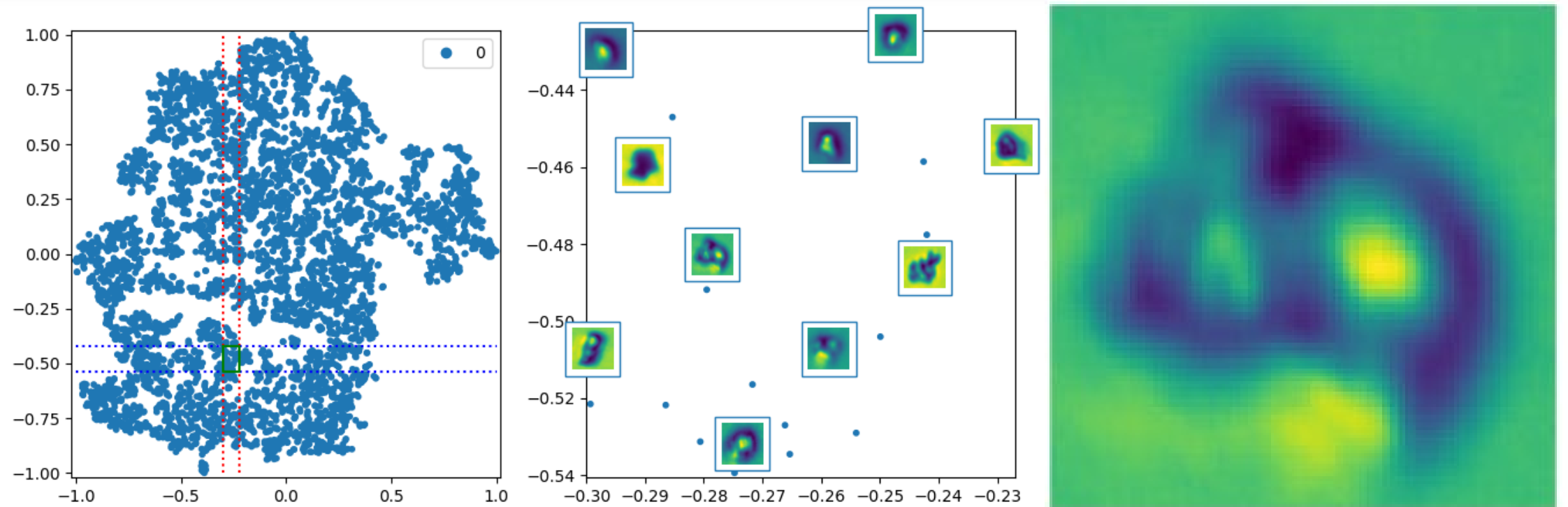
VAE, planar flow Area 2



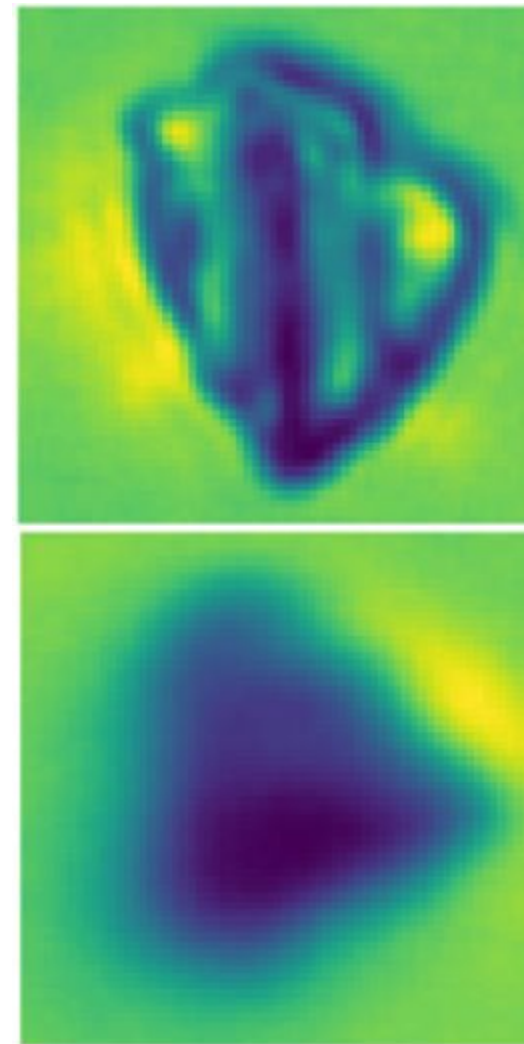
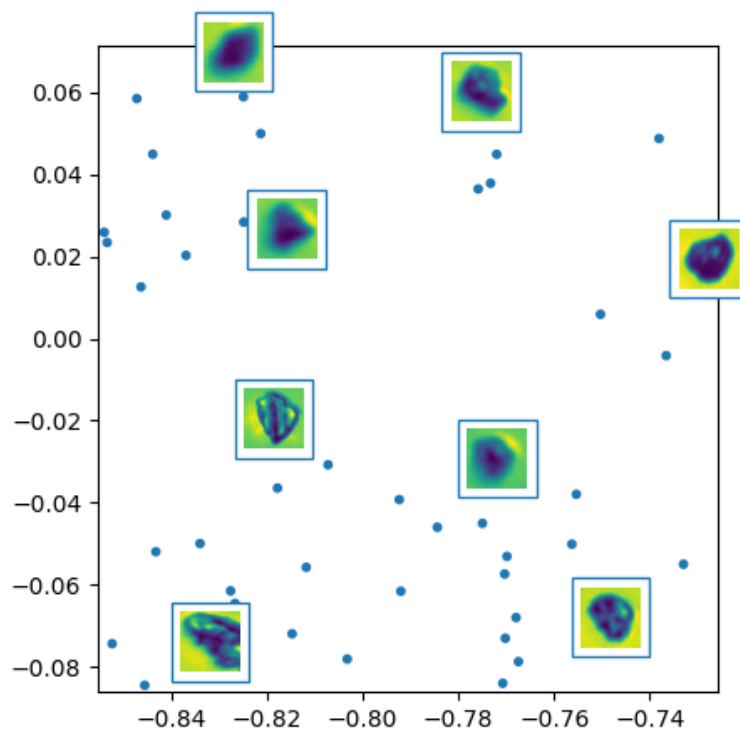
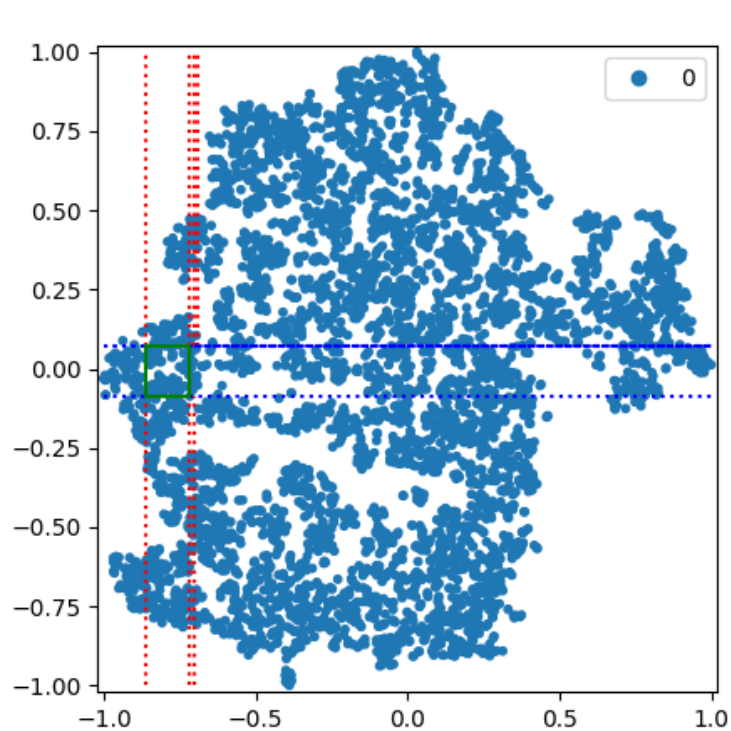
Between area 1 and 3



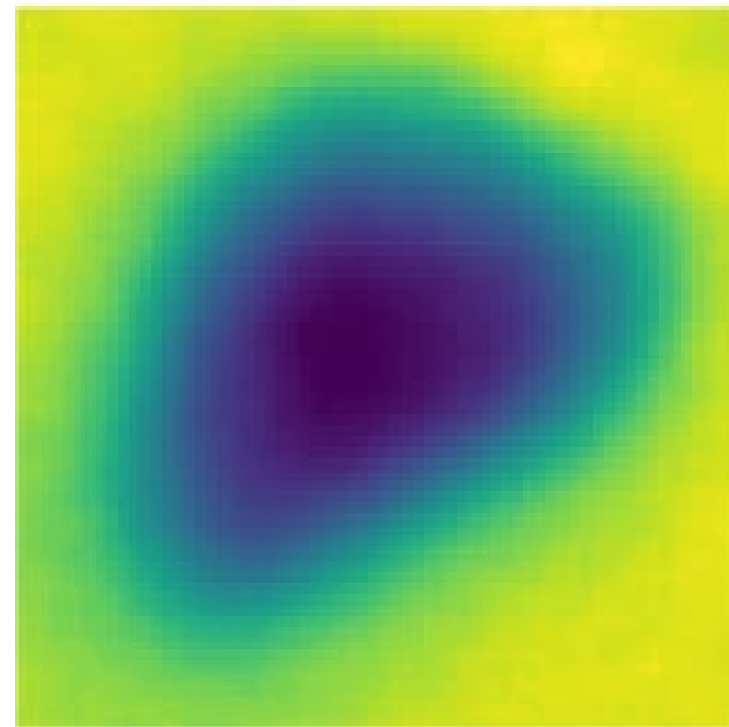
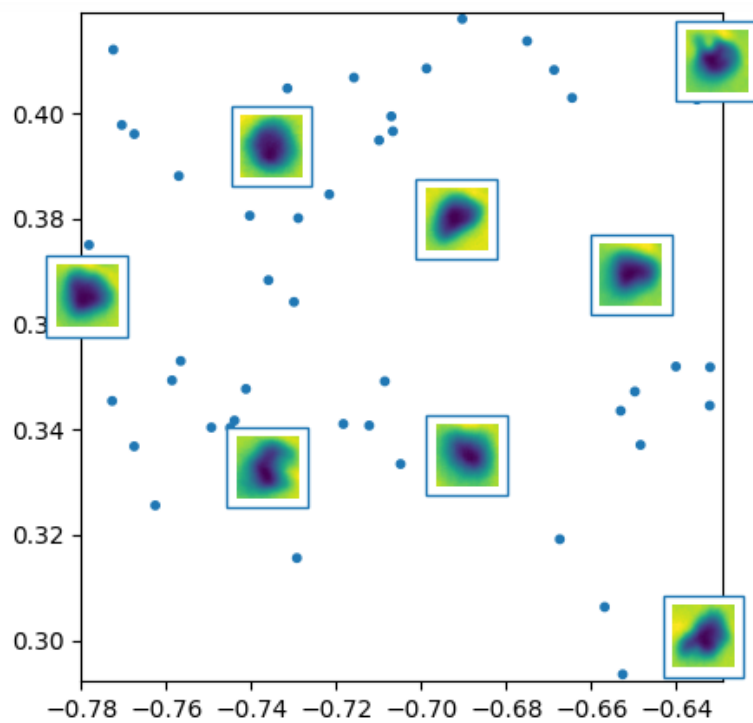
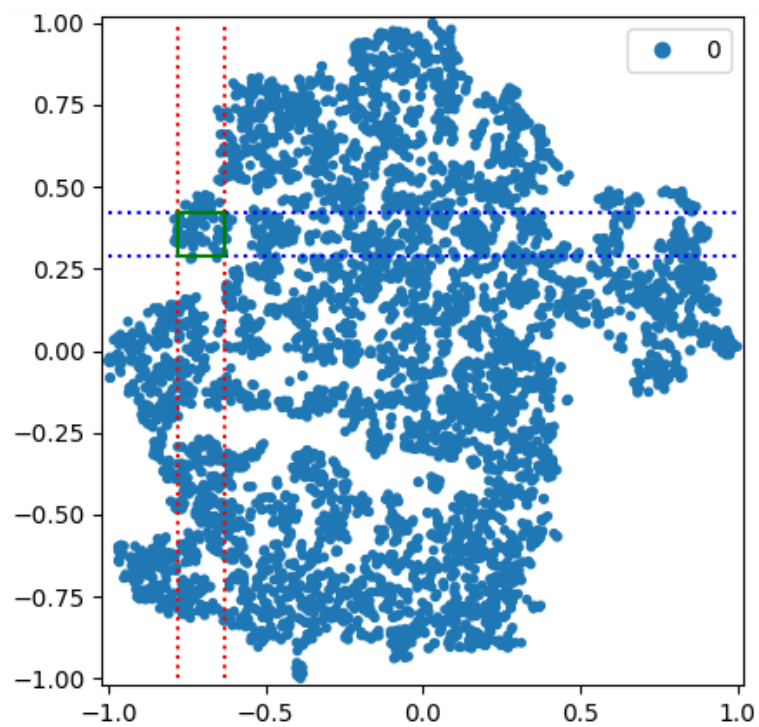
VAE, planar flow Area 4



VAE, planar flow Area 5



VAE, planar flow Area 6



Peruvian Dataset with Labels obtained through the CNN-Classifier

