

# Supporting the analysis of a large coin hoard with AI-based methods

Chrisowalandis Deligio\*<sup>1</sup>, Karsten Tolle<sup>1</sup>, David Wigg-Wolf<sup>2</sup>

<sup>1</sup> Goethe Universität Frankfurt am Main - Germany

<sup>2</sup> Römisch-Germanische Kommission – Germany

\*Corresponding author

Correspondence: Deligio@em.uni-frankfurt.de

## ABSTRACT

In our project "Classifications and Representations for Networks: From types and characteristics to linked open data for Celtic coinages" (ClareNet) we had image data for one of the largest Celtic coin hoards ever found: Le Câtillon II with nearly 70,000 coins. In the initial stages of our approach, the main problem was how to deal with the dataset without having any information about it. First, we separated the dataset into groups of coins of different sizes using object recognition combined with the scale contained in the images. The main approach was to treat the coins independently of the underlying classification and analyse how an unsupervised method could group them. We later evaluated our results against the table provided and produced by the expert team. In addition, we have reviewed these expert classifications to improve them and provide a quality check, but also to get a better understanding of how the experts classify the coins, especially for those in poor condition. Additionally, we took a closer look at a single class and tried to identify coins struck with the different dies used. The phases of our work have been presented at CAA 2022 in Oxford and at CAA 2023 in Amsterdam.

**Keywords:** machine learning, celtic coins, unsupervised learning, classification, object detection

## Introduction

30 In our Classifications and Representations for Networks (ClaReNet) project, we are exploring and  
 31 evaluating computer-based methods applied to three different Celtic coin series. One of them is the stater  
 32 attributed to the Coriosolitae in the hoard of Le Câtillon II found in Jersey in 2012. We are grateful to our  
 33 collaboration partner Philip de Jersey and Jersey Heritage for allowing us to work on this huge dataset of  
 34 120,000 images (around 60,000 per coin face). Before the work began, Jersey Heritage invested a huge  
 35 amount of labour and time (including 25 volunteers) to dismantle the hoard, take the photos and make a  
 36 first identification of each coin. They also provided us with this data during the project. Our goal was to  
 37 support the numismatic process, but also to show the potential of machine learning based methods. We  
 38 aimed to evaluate different issues, from pre-sorting and classifying to the recognition of different dies. The  
 39 main focus was on the use of unsupervised methods to also present how to approach an as yet unknown  
 40 dataset, i.e. without any information other than the images themselves. Finally, we compared the results  
 41 with the classification created by the Jersey team and also involved the expert into our process. This paper  
 42 also presents tools, visualisations and extensions that have proven useful for the task, in communicating  
 43 with the numismatists and integrating their opinion.

44

## Data

45 The data set made available to us includes about 60,000 photos per coin face and contains various coins  
 46 including staters, quarter staters and petit billions. In addition to the coins, the photos mostly include a  
 47 scale and further information such as the assigned identification number (ID) or occasionally the assigned  
 48 class. Both the coins and the photos are of mixed quality: there are broken, worn and corroded coins, as  
 49 well as blurred, overexposed and underexposed photos (fig. 1). We focused on the staters, as they  
 50 represent the largest part of the hoard (about 50,000) and their research is also more advanced. The staters  
 51 are generally divided into six classes (I to VI, Appendix 1), which were originally formed on the basis of the  
 52 obverse, or 'heads' face (Colbert de Beaulieu 1957), therefore we concentrated our work on the obverse  
 53 too. The information we had to evaluate the process is:

- 54 • Staters have an average diameter of about 22 mm, quarter staters and petit billions about 13  
 55 mm.
- 56 • During the project, we received a table from the expert with information about the kind of  
 57 coin (stater, quarter stater, ...) and the assigned class.
- 58 • For one class of the staters (VI ~ 1300 images) there is an unpublished die study.
- 59



60

61

Figure 1 - Variation of the photos and conditions of the coins. (Photos: Jersey Heritage)

62

## Overview of the pipeline

63 With the intention of extracting and analysing the stater, we followed a divide and conquer  
64 methodology. It aims at sequentially dividing the dataset into smaller batches in order to analyse each one  
65 more efficiently. Our pipeline thus includes the following steps:

66

67 **1. Object Detection** - For the rest of the process, the focus should be on the coin image, therefore  
68 the images need to be cropped. At the same time the size can be calculated by detecting of the  
69 scale, which allows a first sorting.

70

71 **2. Unsupervised learning** - In order to use only the images as input, methods were used that do  
72 not require any further domain knowledge. This step of the pipeline was repeated to extract  
73 groups of high similarity and with the goal of identifying groups similar to the expert's  
74 classification, while removing corroded and worn coins to eliminate the bias they can cause.

75

76 From here, the pipeline splits into two paths.

77

78 **3a. Supervised learning** - The result of step two was checked against the classification of the expert  
79 team and a classification model was trained to reassign the coins from the other batches. The  
80 domain expert was also involved in the process.

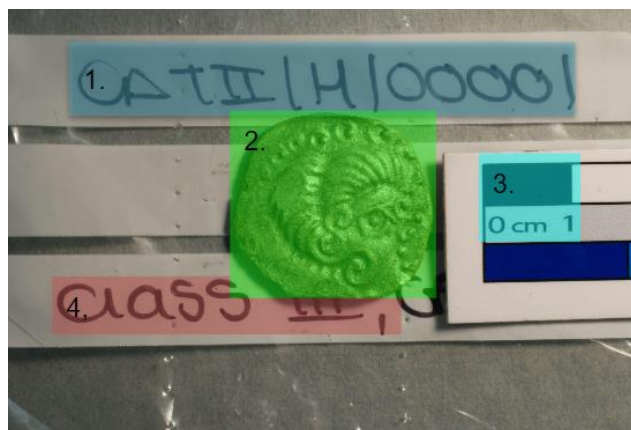
81

82 **3b.** In addition to sorting by class, methods for recognising and sorting by different matrices were  
83 tested. In this case, only the images of one class were used.

84

## Object Detection

85 As already mentioned, the quality of the photos varies, but the format of the photos was consistent.  
86 For most of the photos four areas could be defined, although the last two may be missing: the assigned ID,  
87 the coin, the scale and the class (rare). For us, the coin and the scale were relevant, while the ID can be  
88 retrieved via the file name and the class was available in the spreadsheet provided by the expert and could  
89 be found via the ID (fig. 2).

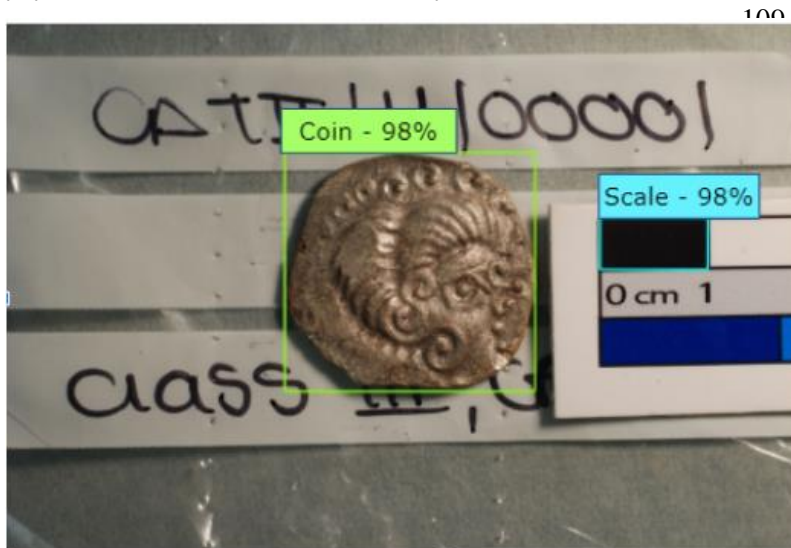


90

91

**Figure 2** - The positions of the defined areas are consistent overall. (Photo: Jersey Heritage)

92 The implementation was done in a supervised process in which two classes, coin and scale, were  
93 defined, and a training and a test dataset were created. The position of the two objects on the images was  
94 consistent and the scale was also a less complex target due to its representation as a simple black box. This  
95 led us to select a relatively small data set in relation to the total number. 100 images were chosen as  
96 training data and 25 as test data; the images were chosen on a more random basis, but it was also  
97 important to cover observations such as small and large coins as well as broken coins. Annotation of the  
98 data was done using the open source tool labelling (Tzutalin 2015). The evaluation of the test dataset after  
99 training gives a mean average precision of 95%. The evaluation for the whole dataset could not be given as  
100 a percentage, but by calculating the size of the coins it was possible to identify outliers and thus improve  
101 the procedure in a targeted way. In addition to the outliers detected by size, we manually re-measured 10  
102 coins to get a feeling for the quality of the result, which turned out to be very accurate. Finally, we cropped  
103 the images and verified whether the coin is still completely displayed. Although 100 coins was only a small  
104 sample, it did prove that this number of images was sufficient for the task. An overall evaluation would  
105 require the annotation of all data, which is time-consuming and not essential for the further process. For  
106 the implementation, we used Tensorflow's Object Detection API<sup>1</sup> and its Model Zoo<sup>2</sup> in order to select a  
107 model architecture. We have not attempted to evaluate different architectures against each other in this  
108 paper. However, for our task we finally decided to use the CenterNet Hourglass104 512x512 by Duan et al.



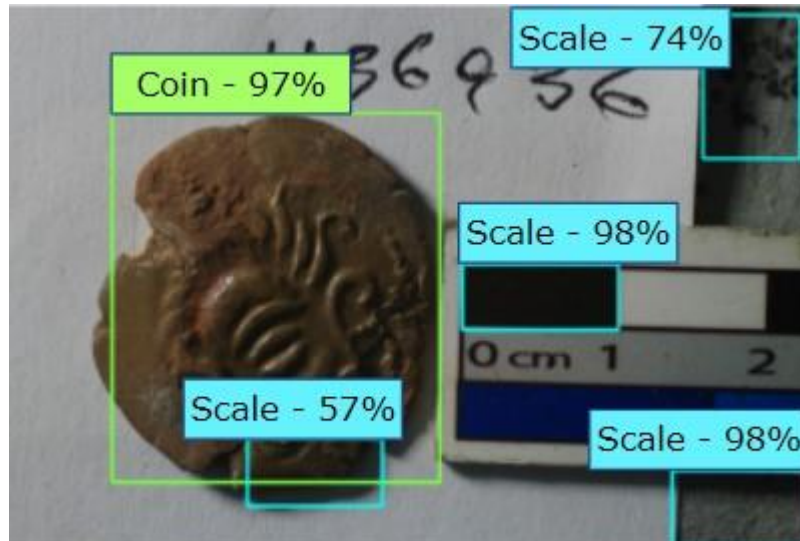
**Figure 3** - Optimal Prediction of the model. Calculated values: height: 2.321cm, width: 2.194cm. (Photo: Jersey Heritage. Graphic: C. Deligio, Big Data Lab)

129 maximum diameter. Therefore, we decided to take only two measurements: the width and height of the  
130 detected coin bounding box, dividing these two values by the width of the detected bounding box of the  
131 scale, which gives the value in centimetres. The cases where the size differed greatly from the others, and  
132 was therefore to be classified as an outlier, were considered separately. Figure 4 shows such a case. There  
133 are several reasons that could lead to an incorrect calculation. It can be seen that darker areas are classified  
134 as scale, which is not surprising when the target is a black box. In figure 4 there is also an area that achieved  
135 the same percentage as the scale itself, which can cause selection problems during the calculation. To  
136 counteract such cases, these and other suitable images (e.g. images with shadows) have been annotated,  
137 thus broadening the training base.

<sup>1</sup> [https://github.com/tensorflow/models/tree/master/research/object\\_detection](https://github.com/tensorflow/models/tree/master/research/object_detection)

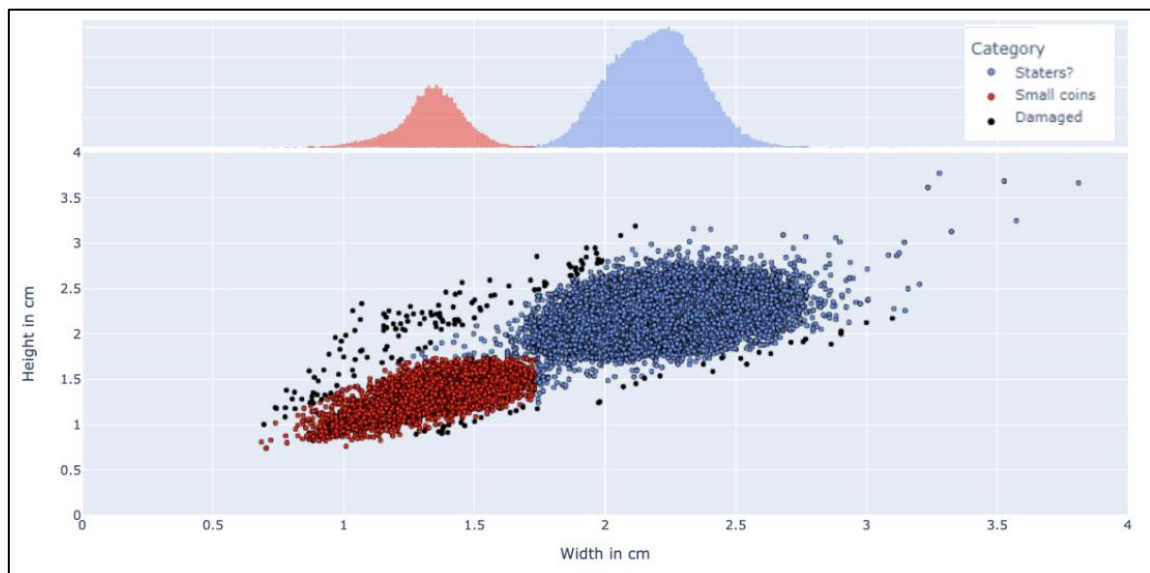
<sup>2</sup> [https://github.com/tensorflow/models/blob/master/research/object\\_detection/g3doc/tf2\\_detection\\_zoo.md](https://github.com/tensorflow/models/blob/master/research/object_detection/g3doc/tf2_detection_zoo.md)

<sup>3</sup> <https://cocodataset.org/#home>



138  
139  
140 **Figure 4** - Shadowy areas can lead to a wrong prediction by the model, resulting in an incorrect size calculation. (Photo: Jersey Heritage. Graphic: C. Deligio, Big Data Lab)

141 The calculation was made for each photo, i.e. for both sides of the coin, resulting in four measurements  
142 for each coin. Comparing these also provides confirmation, as cases where the values are very different  
143 were examined more closely. The calculated values in centimetres can be viewed in a scatter plot in Figure  
144 4. The colour coding is as follows: Coins with a height and width deviation of more than 40% are defined  
145 as damaged and marked in black. Small coins (probably quarter stators and petit billions) are marked red  
146 and large coins (probably stators) are marked blue. Where we can separate the different groups exactly  
147 needs to be examined more closely, but the visualisation shows that there is a gap at about 1.75cm, so we  
148 have chosen this as the separation. The blue area is the focus of our work because it should contain the  
149 stators, but as we want to analyse it in more detail first, we will call this group the 'Stators?'.



150 **Figure 5** - Scatter plot of the approximated diameter. (Graphic: C. Deligio, Big Data Lab)

151 If we look at the peaks of the two point clouds, most of the elements are at around 1.3cm and 2.2cm.  
152 These two values correspond to the information provided by the expert and thus provide the first  
153 confirmation of our process. The dataset can therefore be divided into four groups: "Stators?" (54.227 -  
154 91.29%), "Small Coins" (3.340, 5.64%), "Damaged" (97, 0.16%) and "Not Detected" (1.778, 2.91%). The  
155 latter contains photos where no scale was available or it was not detected, and therefore the calculation  
156 could not take place.

157

## Unsupervised Learning

158

159

160

161

162

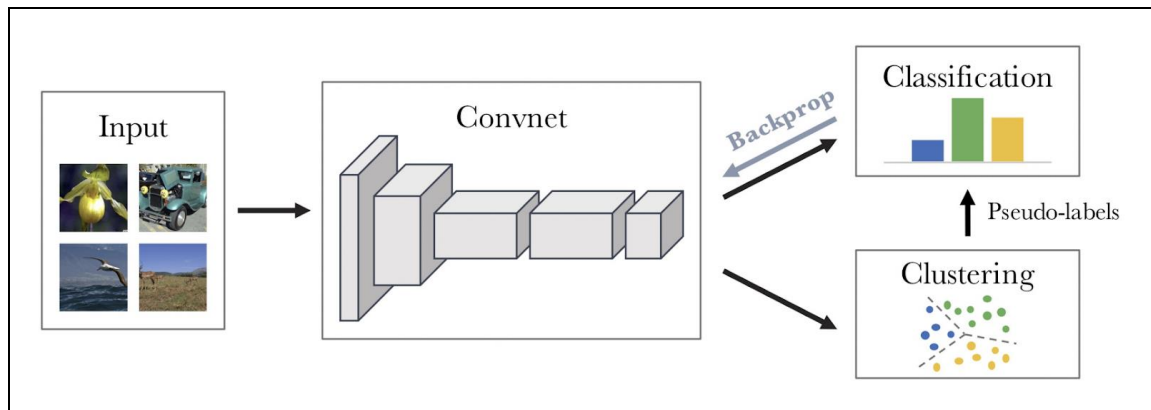
163

164

165

166

Taking our divide and conquer approach one step further, we selected the "Stater" group as the dataset for unsupervised learning. The main aim was to see if we can achieve a pre-sorting and to see how far it complies to the expert classification. Deep Learning and convolutional neural networks (CNNs) were the first choice for image classification, and a promising approach we found was the DeepCluster (Figure 6) method developed by Caron et al. (2018). This combines a convolutional neural network, which was what we wanted to use, with a clustering algorithm to train a CNN in an unsupervised manner. The idea in this approach is to use the generated clusters as pseudo-labels to train the CNN, and the extracted features in turn serve as input to the clustering algorithm. This process is then repeated for the desired number of epochs.



167

168

169

**Figure 6** - DeepCluster implements a method of training a CNN in an unsupervised way. (Caron et al. 2018)

170

171

172

173

174

175

176

177

178

179

180

In the Paper by Caron et al. 2018, the method was used with the two CNN architectures alexnet and VGG16. We used the architecture with VGG16, and as the clustering algorithm k-Means. Thus, the required inputs were the images and the number of desired clusters (k). The choice of k was initially a challenge, for the paper recommends a much larger k (e.g. the best results are obtained with a k 10 times larger than the actual number of classes). On the other hand, the number of expected classes is known (6), but since we wanted to analyse the dataset without using any information, we started with a k equal to 100. To measure another factor of the effectiveness of the method, for our first exercise we entered both the obverse and reverse photos to see if they would be separated. For the evaluation of the resulting clusters, we avoided the use of additional information and performed it manually. Based on this manual evaluation, the following observations were made (fig. 7):

181

182

183

184

185

186

187

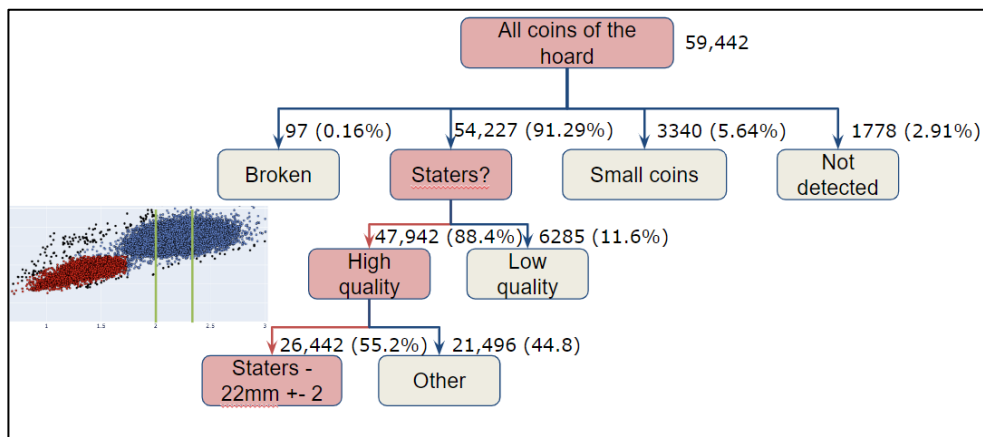
- Obverse and reverse were generally not mixed within the clusters → showing that the method is already working at a high level.
- Coins in poor condition were grouped together → showing the potential to clean up the dataset for further processing.
- Clusters with different levels of wear were identified.
- There were mixed clusters with no common features → CNNs are complex and have been described as black boxes, so sometimes it is not clear how the results were obtained.



188  
189 **Figure 7** - Example of clusters generated. (Photos: Jersey Heritage)

190 Many clusters emerged that showed a strong similarity (based on our manual evaluation). Clusters with  
191 corroded and poorly preserved coins which had been identified were sorted out. This allowed us to divide  
192 the dataset into "High Quality" - clusters with a high similarity and well preserved coins - and "Low Quality"  
193 - corroded and worn coins - or also clusters with a significant mix. In order to provide another degree of  
194 certainty, we focused on the coins at the centre of the size point cloud, i.e. 2.2cm +/- 0.2cm (fig. 5). In this  
195 way, we wanted to ensure that only staters were present in the dataset (fig. 8). Also, we only used the  
196 obverse images.

197



198 **Figure 8** - Using the divide and conquer methodology, the data set could be divided step by step into  
199 more easily analysable parts. (Graphic: C. Deligio, Big Data Lab)

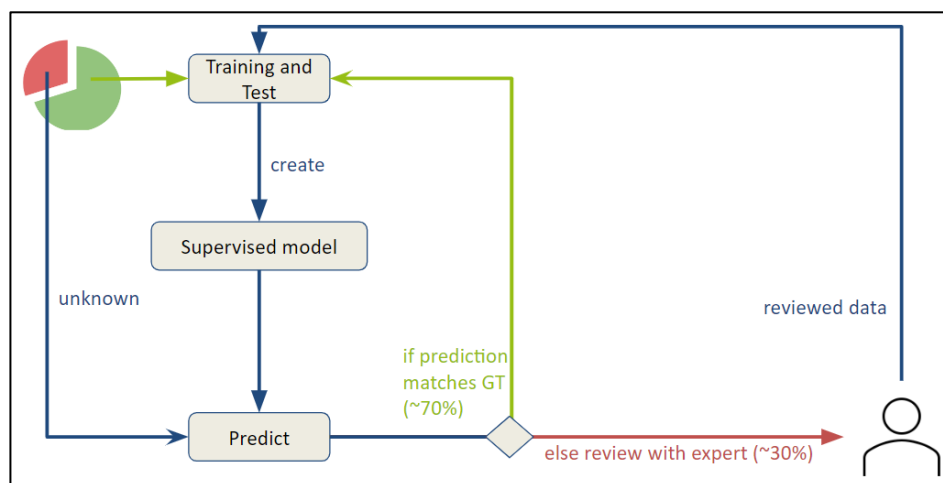
200 This new, reduced dataset (ca. 26,000 images) was now to be divided into clusters (25) once again, and  
201 when checking the results we evaluated them against the spreadsheet of the expert team. The first  
202 evaluation was to determine whether our selected dataset really did contain only staters. This could be  
203 confirmed. By comparing our results with the experts' spreadsheet, we could also calculate exactly which  
204 classes each cluster was composed of: the result was that 18 out of 25 clusters contained at least 70% of  
205 elements of only one single assigned class (15,063 images). Within this threshold, only 8% do not  
206 correspond to an assigned class (1208 images). The evaluation of the cluster is presented in Appendix 2.  
207 One drawback was that we did not manage to find any clusters of class VI coins which, with about 1300  
208 occurrences, is by far the smallest class in the dataset. The images of class VI coins were mostly mixed into

209 clusters of class V or sometimes IV, which is probably a result of the similarity of the classes. In total, 13,855  
 210 coins (23% of the total data set) were confirmed by the comparison with the spreadsheet This data set,  
 211 verified by two systems, by deep learning and the team's classification, now formed the basis for the  
 212 supervised approach and verification of the coins previously excluded.

213 **From unsupervised to supervised**

214 The reduced and validated dataset was used as the basis for building a supervised, trained CNN model.  
 215 As the dataset was highly unbalanced (ranging from 615 to 5317 images per class) and we could not  
 216 automatically extract class VI with our method, we adapted the dataset slightly. As we did not want to lose  
 217 class VI, we added the coins validated by the expert. As the unbalance was high, we rebalanced the dataset  
 218 by downsampling to the smallest class. The goal was to increase the training material step by step. As a  
 219 choice of CNN architecture, we started with VGG16. With the supervised model, we intended to evaluate  
 220 and revise the other batches in our tree (fig.8). At the same time, we involved the domain expert, because  
 221 the predictions had two outcomes: if the prediction matches the assigned one, it is basically confirmed.  
 222 However, cases where the prediction differs are not automatically classified as false, but instead saved  
 223 separately, to be reviewed. Confirmed coins were added to the training set so that the model constantly  
 224 receives more material (fig. 9). This means we improved the model and checked the data quality at the  
 225 same time.

226



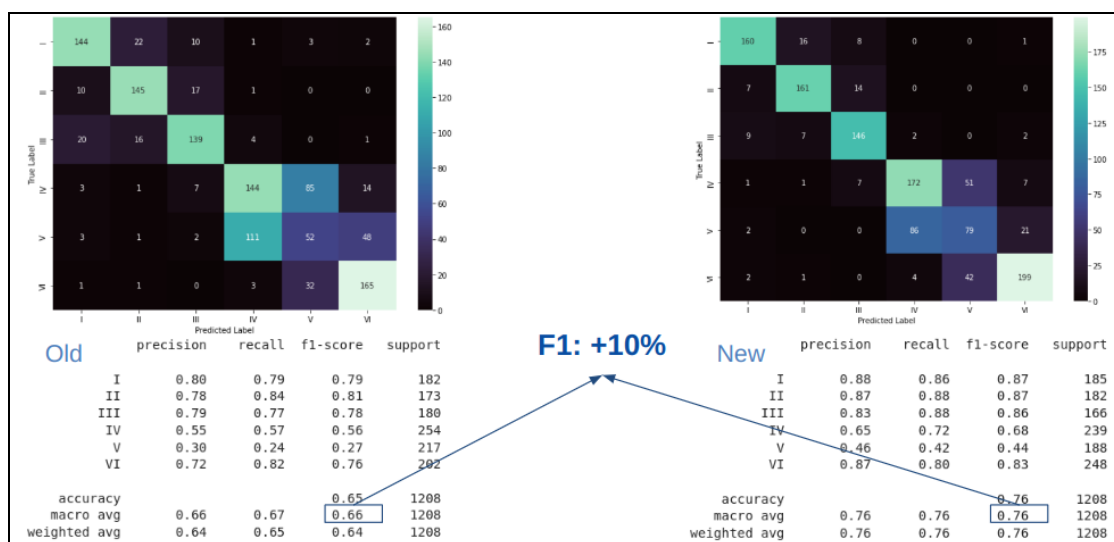
227 **Figure 9** - The 13,855 images from the previous step were selected as the first training base (green).  
 228 The 1208 wrongly selected ones as the first test set (red). The percentages are the results of the  
 229 predictions on the test set. (Graphic: C. Deligio, Big Data Lab)  
 230

231 Whether the prediction or the given class was correct had to be decided by someone with domain  
 232 knowledge, and therefore we involved the domain expert. We also used this step to design a kind of  
 233 experiment. We created a list of coin images with two options, prediction and original assignment, but  
 234 which were masked so that the expert did not know which was which, along with the ID. The domain expert  
 235 could then choose one of the options or specify a new one. Comments were also welcome as especially in  
 236 the case of difficult decisions, they helped us to understand them as non-experts. In our first test set  
 237 (explained in Figure 9), the model was off by 30% (328 out of 1208 images). The review by the expert for  
 238 these cases could be divided into four cases:

- 239
- 240 1. The assigned class was actually wrong and was improved by the model --> data quality
  - 241 improvement (115 cases - ~35%).
  - 242 2. The determination was not clear --> problematic cases (26 cases - ~8%)
  - 243 3. The model was wrong (175 cases - ~53% - mostly between class IV and V - 126 cases)
  - 244 4. Both were wrong (12 cases - ~4%)
  - 245



246 The model mainly had problems distinguishing between classes IV and V, but this phenomenon also  
 247 occurs with humans. We also distributed the list to our team (i.e. to non-experts) and asked them to fill it  
 248 out. With the result that exactly these two classes also led to problems. In figure 10 we can see the  
 249 evaluation of the 1208 images with the old classification (the one originally supplied) and with the expert's  
 250 revised classification for these cases. The F1 metric on the revised classification increased by 10%. This  
 251 shows above all that the performance of such a model cannot be calculated exclusively on the basis of  
 252 metrics if the underlying data quality is not given.

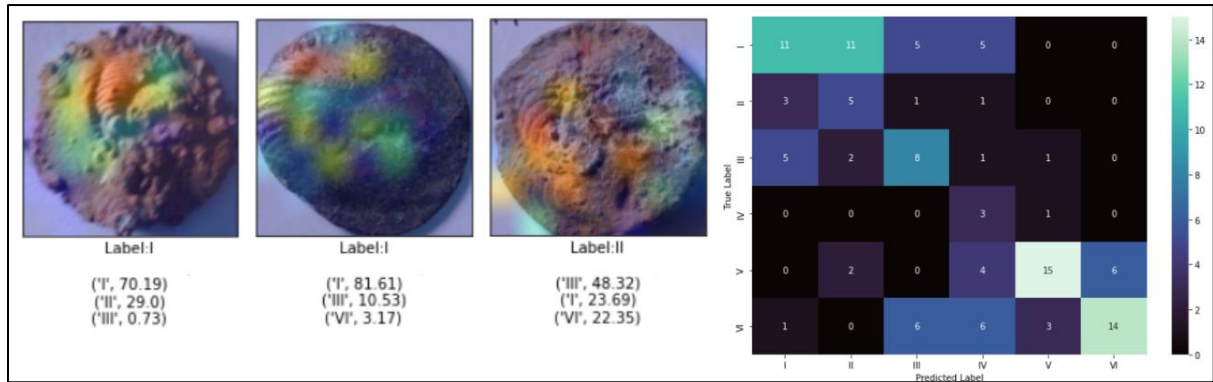


253  
 254 **Figure 10** - Same predictions, different results. Comparison between two classifications (old vs  
 255 revised). (Graphic: C. Deligio, Big Data Lab)

256 This process could now be repeated step by step for the remaining batches that had previously been  
 257 sorted out. Ultimately, this could be used to improve the model and provide higher data quality.

258 An important question we also asked ourselves was how to deal with the coins defined by us as 'low  
 259 quality', or to what extent such a model can help us. Using a random sample of 20 images per class from  
 260 this set, we evaluated a small case study. It was of particular interest to see if the model that had been  
 261 trained on very good images can be applied here. Our case study achieves an accuracy of 47%. Figure 11  
 262 shows the confusion matrix. The figure also visualises three examples with GradCam. It can be seen that  
 263 correct regions (such as the hair or eye) are detected, but that there is also a bias due to the condition. The  
 264 two images on the left are correctly classified by the model. Comparing the right-hand image with the  
 265 images in Appendix 1, it could indeed very well be class III (based on the style of the eye). This was only a  
 266 preliminary test, but it showed the importance of the ground truth, and that it can and should be  
 267 questioned. It is also clear that in order to improve the model, more such material should be integrated to  
 268 counteract the bias of the condition of the coins.

269  
 270  
 271  
 272



273

274

275

276

**Figure 11** - On the left is the visualisation of the prediction with the top 3, on the right is the matrix of the 120 predictions. An F1 value of 44% and an accuracy of 47% were achieved. (Photos: Jersey Heritage. Graphic: C. Deligio, Big Data Lab)

277

### Visualisations and Augmentations



**Figure 12** – Cutout random parts of the image. (Photos: Jersey Heritage. Graphic: C. Deligio, Big Data Lab)

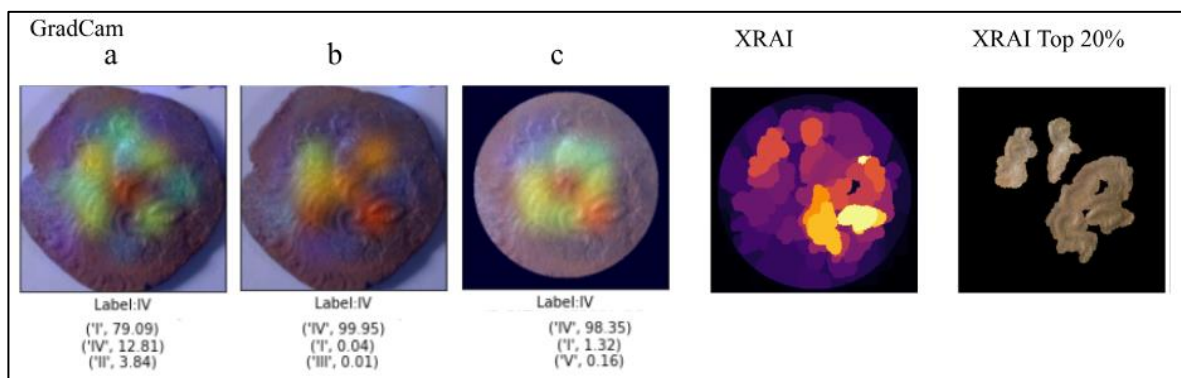
The involvement of the expert also led to exciting insights, e.g. that the nose area plays a central role for him when assigning the class, but the visualisations of the predictions by GradCam, for example, did not reflect this focus. In the example of the coin in figure 13 (left), the right class is indeed in the top 3, but with only a 12.8% certainty. The important features, not always receive as much weight as desired. To address this issue and to try to incorporate the insights of the domain expert, as well as to influence the training process, we tried several augmentation methods. Two of them turned out to be particularly helpful for our case. The cutout method involves hiding parts of the image, thus leading it to pay attention to several other areas (fig. 12). The cutout can be targeted or randomised and used with a fixed seed to replicate the training. Looking at the same coin with the model trained with the cutout augmentation (fig. 13b), we can see that certain areas have a stronger weighting, especially the eye region.

291

292

293

294



295

296

297

298

299

**Figure 13** - GradCam visualisation and top 3 prediction of a model trained based on full coin (a), cutout (b), circle crop (c) images. Colour scale: blue (weak) - red (strong). On the right a visualisation with the use of XRAI (circle crop images) and cropping the top 20% area. (Photos: Jersey Heritage. Graphic: C. Deligio, Big Data Lab)

300

301

Another augmentation, we simply call it the circle crop, was chosen because sometimes the edge of the coin, which can be very irregular, is focussed on and so can also cause noise. To counteract this, we

302 applied a simple circle crop oriented on the centre of the image to remove the edges. Figure 13 (c) shows  
303 that the focus of GradCam was very much in the centre, but less weighted compared to the cutout  
304 augmentation. In both cases you can also see that the class has been correctly identified.

305  
306 As one can see, it is possible to direct the focus a little, and this offers the possibility of incorporating  
307 domain knowledge and to some extent also preferences. When it comes to explaining how a CNN works  
308 to a numismatic team, things can quickly get complicated. CNNs are difficult to understand due to their  
309 complexity and large number of parameters. To overcome this, there are various methods (SHAP, Lundberg  
310 and Lee 2017; LIME, Riberio et al. 2016; XRAI, Kapishnikov et al. 2019; and many more) of explanation,  
311 including visual ones, such as the GradCam method shown here. We recommend trying different methods  
312 and to communicate with the team to find a suitable one. In addition to GradCam, we also decided to use  
313 XRAI. While GradCam expands from one point and is more coherent, XRAI is meant to be independent and  
314 more focused on the relevant features regardless of the location (Kapishnikov et al. 2019). The  
315 implementation<sup>4</sup> used also offered the option to extract the most important features instead of displaying  
316 a heat map, which was well received by the team and could be interpreted quickly

### 317 **Implementing a die study**

318 For a die study we dug further into a single class in order to try to distinguish individual coin dies within  
319 it. As noted above, for one of the six classes (VI with about 1300 images) an unpublished die study was  
320 available to us in order to evaluate our applied methods. The recognition of dies brings another level of  
321 challenge for, although we are dealing with less data, the coins are very similar (same class) and there are  
322 many dies. For example, previously we had about 60,000 images with the goal of distinguishing six classes,  
323 but now we had only 1300 images with over 30 defined dies (based on the unpublished die study). For the  
324 implementation, we tested three methods against each other:

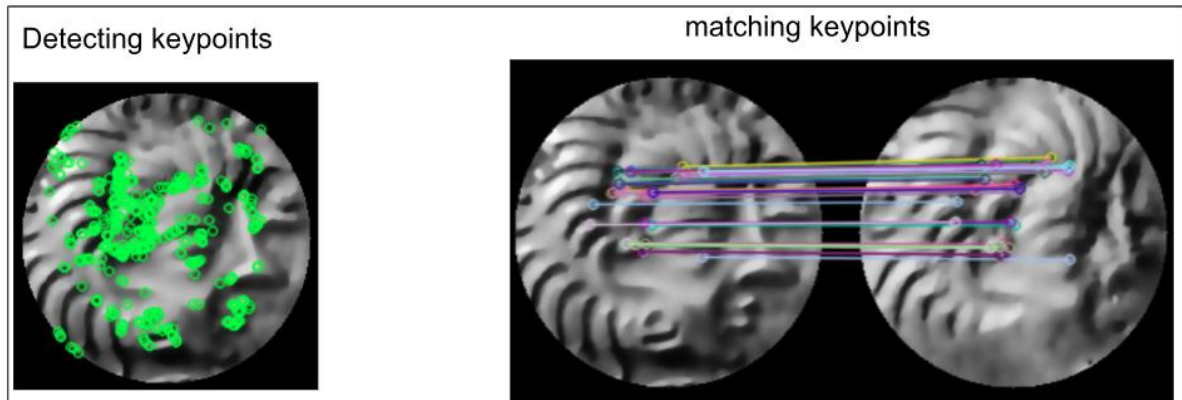
- 325
- 326 1. Reapplying DeepCluster.
- 327 2. Using our trained supervised model to extract features and then cluster them.
- 328 3. Algorithms comparing the key points in the image, which has also been successfully used in  
329 other publications.
- 330

331 The first two methods are similar in principle but differ in the trained CNN, in (1) it is trained from  
332 scratch using the DeepCluster algorithm, in (2) we used our trained model on the six classes and used it to  
333 extract features. (3) is a very different method to the ones used so far and will be discussed briefly here.  
334 The algorithms used are from the field of image matching. The best known are probably SIFT (Scale-  
335 invariant feature transform by Lowe, 2004), and SURF (Speeded Up Robust Features by Bay et al., 2006).  
336 As the first two are patented algorithms, a popular open source alternative is ORB (Oriented FAST and  
337 rotated BRIEF by Rublee et al., 2011). Such algorithms have also been successfully applied to ancient coins  
338 in various publications (Kampel and Zaharieva 2008; Taylor 2020; Heinecke et al. 2021). For our procedure  
339 we used ORB. There are also some pre-processing steps that had a positive effect on the results in terms  
340 of reducing bias (like scratches from the usage of the coin). The images were converted to greyscale in  
341 order to avoid a colour bias. The images were also blurred and contrast adjusted (see Heinecke et al. 2021).  
342 Finally, a circle crop was applied to remove the edges of the coins and to focus only on the motif (fig. 14  
343 shows the output).

344

---

<sup>4</sup> <https://github.com/PAIR-code/saliency>



345

346  
347  
348

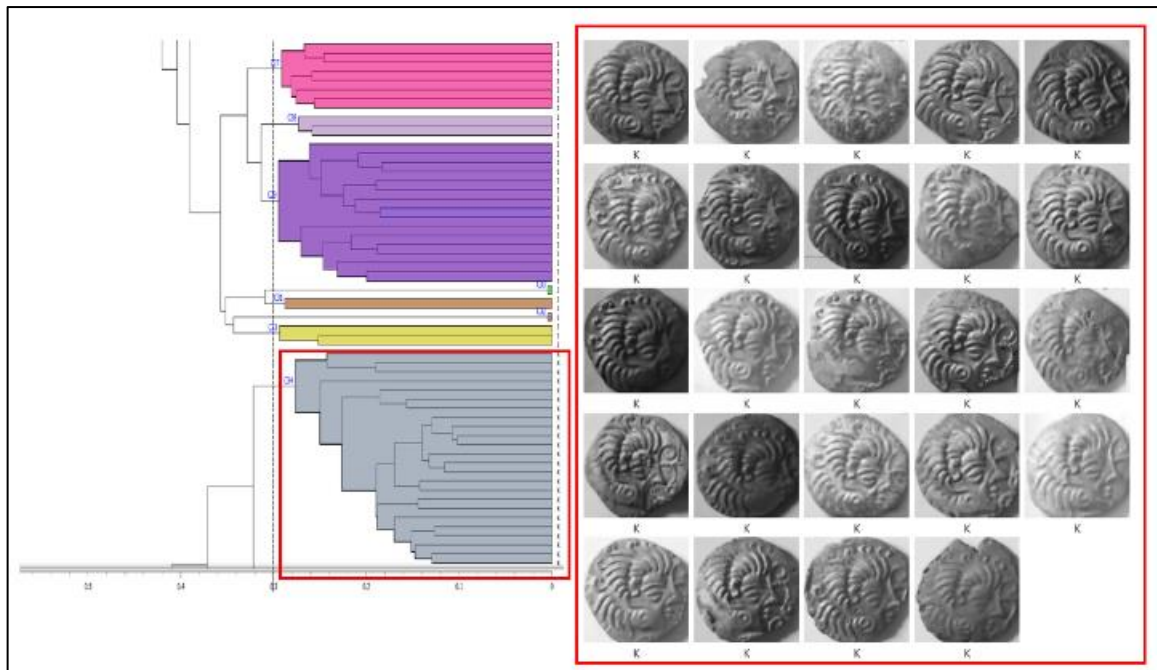
**Figure 14** - On the left, an image with detected keypoints. On the right, an example of two supposedly identical die pieces (according to the GT) and their matches. (Photos: Jersey Heritage. Graphic: C. Deligio, Big Data Lab)

349 The process is then as follows, the first step consists of the key point detection algorithm and matching  
350 the key points between two images. This comparison is carried out in pairs between all images and the  
351 matches found are captured as a vector for each image (resulting in an  $n \times n$  matrix). For the exact method  
352 of calculating the key points, we recommend the ORB publication (Rublee et al., 2011).

353

354 The second step is the same for all three methods, the features from (1) and (2), which are also stored  
355 as a vector, and the result from (3) are all used as input to a clustering algorithm, in this case we used  
356 hierarchical clustering. We used the [Orange Data Mining](#) tool for calculating distances (based on Spearman  
357 distance metric), clustering and visualisation. Figure 15 shows the visualisation of the clustering as a  
358 dendrogram, in which we compare the result of the clustering with the existing die study of the expert.  
359 In order to evaluate the methods equally, there are now various possibilities; we decided to evaluate them  
360 all with the same distance value. We started with the image matching method (3) and the value 0.3  
361 proved to be optimal and was also chosen for the other two methods for direct comparison. Table  
362 1 is a summary of the results that were obtained with this value. It shows the best value achieved within a  
363 cluster in terms of the number of coins from an individual die identified, the mean value of all clusters and  
364 the total number of clusters formed.

365



366

367  
368

**Figure 15** - On the left, part of the dendrogram. On the right, an overview of a cluster successfully containing the coins of one die. (Photos: Jersey Heritage. Graphic: C. Deligio, Big Data Lab)

369  
370  
371  
372  
373  
374  
375

It is noticeable that the DeepCluster method performs better with a larger amount of data and in this case struggled with the relatively high number of classes on this small amount of images. The second method, our supervised trained model, performed slightly better. The advantage was that the model already existed and no additional training time was required. The third method based on ORB worked best. Looking more closely at the results of the third method, of the 256 clusters, 208 clusters had at least 70% coins from the same die. Furthermore, 194 of them had only coins of exactly one die. The 194 cluster cover 489 coins and that was 40% of the dataset considered here.

376  
377

**Table 1** - To compare the methods we used the same threshold (0.3), we calculated two values to get an impression of the performance: Highest matching cluster with gt and the mean over all clusters.

<b>Method</b>	<b>DeepCluster (k=15)</b>	<b>Supervised model (CNN)</b>	<b>Keypoint detection &amp; matching (with ORB)</b>
<b>Nr. of clusters at distance threshold (0.3)</b>	45	172	256
<b>Highest correspondance</b>	60%	75%	100%
<b>Mean</b>	24%	37%	84%

378

379

### **Recapitulation and outlook**

380  
381  
382  
383  
384  
385  
386  
387  
388

We started this research by treating the dataset as a case study of a new find, with no information available at the outset. With the first step of object detection, it was possible to automatically crop the images and, with the scale, calculate the size of the coins and perform some sorting, which helped to identify the staters. The next step of unsupervised learning was to see how far the dataset could be sorted or grouped, and whether it could even match the classification of the expert. By repeatedly applying the method with manual evaluation and then combining it with the size calculation, it was possible to extract a dataset containing only staters and to extract about 14,000 coins, which were correctly grouped according to the provided spreadsheet.

389  
390  
391  
392

Re-evaluating the data with a supervised method leads to an improvement in data quality and in critical cases the AI's and the expert's decision can be compared . It is necessary to involve the domain expert in the process.

393  
394  
395

Specific additions can influence the choice of features and bring a model closer to the expert's expectations.

396  
397  
398  
399  
400

Sometimes it is also useful to use have a look on other algorithms (in terms of complexity and computing power required), in our case image matching (ORB), especially with small datasets (many classes, high similarity of data), as is shown by the positive results of our die study and of tests presented in other publications. Selecting the right approach and algorithm is still a challenge, also for IT experts. This is partly due to the fact that each approach is also accompanied by various subtasks (e.g. different ways for

401 preprocessing or augmentation) and possibilities for fine tuning (e.g. the number of clusters,  
402 hyperparameter settings, choosing the best loss function).  
403  
404 The results show that semi-automatic support can be helpful in conducting sorting, classification, or even  
405 a die study. A tandem approach where domain and IT experts work closely together will probably have the  
406 best success rate.

## 407 **Data, scripts, code, and supplementary information availability**

408 Implementations used in this paper:: [https://github.com/Frankfurt-BigDataLab/2023\\_CAA\\_ClaReNet](https://github.com/Frankfurt-BigDataLab/2023_CAA_ClaReNet)

409 Dataset used in this paper: <https://doi.org/10.34780/kzw0-r608>

410 Official website of the dataset: <https://catalogue.jerseyheritage.org/categories/celtic-coin-hoard/>

411 Official implementation of DeepCluster by Caron et al.: <https://github.com/facebookresearch/deepcluster>

412 Annotation tool used: <https://github.com/heartexlabs/labelImg>

413 Implementing an object detection model: <https://github.com/sglv1adi/TensorFlowObjectDetectionTutorial>

414 For implementing a supervised model: [https://www.tensorflow.org/tutorials/images/transfer\\_learning](https://www.tensorflow.org/tutorials/images/transfer_learning)

415 Useful augmentations library: <https://albumentations.ai/>

416 Tool for visualising results (and more): <https://orangedatamining.com/>

## 417 **Conflict of interest disclosure**

418 The authors declare that they comply with the PCI rule of having no financial conflicts of interest in relation  
419 to the content of the article.

## 420 **Funding**

421 The Classifications and Representations for Networks project is funded by the German Federal Ministry of  
422 Education and Research (BMBF).

## 423 **References**

424 Bay, Herbert, Tinne Tuytelaars, Luc Van Gool. 2006. *SURF: Speeded Up Robust Features*. In: Leonardis, A.,  
425 Bischof, H., Pinz, A. (eds) *Computer Vision – ECCV 2006*. ECCV 2006. Lecture Notes in Computer Science,  
426 vol 3951. Springer, Berlin, Heidelberg. [https://doi.org/10.1007/11744023\\_3](https://doi.org/10.1007/11744023_3)

427  
428 Caron, Mathilde, Piotr Bojanowski, Armand Joulin, Matthijs Douze. 2018. *Deep Clustering for Unsupervised*  
429 *Learning of Visual Features*. <https://doi.org/10.48550/arXiv.1807.0552>

430  
431 Colbert de Beaulieu, Jean-Baptiste. 1957. *Le trésor de Jersey-11 et la numismatique celtique des deux*  
432 *Bretagnes*. *Revue Belge de Numismatique* Vol. 103. 47-88.

433  
434 Duan, Kaiwen, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang and Qi Tian. 2019. *CenterNet: Keypoint*  
435 *Triples for Object Detection*. <https://doi.org/10.48550/arXiv.1611.10012>

436  
437 Heinecke, Andreas, Emanuel Mayer, Abhinav Natarajan, Yoonju Jung, *Unsupervised Statistical Learning for*  
438 *Die Analysis in Ancient Numismatics*. *Computer Vision and Pattern Recognition* 2021.  
439 <https://doi.org/10.48550/arXiv.2112.0290>

440 Huang, Jonathan, Vivek Rathod, Chen Sun, Menglong Zhu, Anoop Korattikara, Alireza Fathi, Ian Fischer,  
441 Zbigniew Wojna, Yang Song, Sergio Guadarrama and Kevin Murphy. 2017. *Speed/accuracy trade-offs*  
442 *for modern convolutional object detectors*. <https://doi.org/10.48550/arXiv.1611.10012>

443





Appendix



479

480  
481

Appendix 1 - The six classes of staters as defined by numismatists. (Photos: Jersey Heritage)

Cluster	Class_I	Class_II	Class_III	Class_IV	Class_V	Class_VI	Other	Total	Class_I_%	Class_II_%	Class_III_%	Class_IV_%	Class_V_%	Class_VI_%
0	0	8	649	9	0	0	11	677	0.0	0.012	<b>0.959</b>	0.013	0.0	0.0
1	0	21	646	2	0	0	6	675	0.0	0.031	<b>0.957</b>	0.003	0.0	0.0
2	1140	24	7	0	7	0	7	1185	<b>0.962</b>	0.02	0.006	0.0	0.006	0.0
3	0	0	0	153	966	55	1	1175	0.0	0.0	0.0	0.13	<b>0.822</b>	0.047
4	22	1120	24	0	0	0	34	1200	0.018	<b>0.933</b>	0.02	0.0	0.0	0.0
5	2	721	7	0	0	0	8	738	0.003	<b>0.977</b>	0.009	0.0	0.0	0.0
6	3	1	2	501	802	130	20	1459	0.002	0.001	0.001	0.343	0.55	0.089
7	0	0	0	6	762	197	0	965	0.0	0.0	0.0	0.006	<b>0.79</b>	0.204
8	995	1016	660	572	1214	418	366	5241	0.19	0.194	0.126	0.109	0.232	0.08
9	0	7	490	27	0	0	2	526	0.0	0.013	<b>0.932</b>	0.051	0.0	0.0
10	13	842	2	0	5	0	18	880	0.015	<b>0.957</b>	0.002	0.0	0.006	0.0

Cluster	Class_I	Class_II	Class_III	Class_IV	Class_V	Class_VI	Other	Total	Class_I_%	Class_II_%	Class_III_%	Class_IV_%	Class_V_%	Class_VI_%
11	647	293	349	29	23	5	34	1380	0.469	0.212	0.253	0.021	0.017	0.004
12	0	11	543	12	0	0	3	569	0.0	0.019	<b>0.954</b>	0.021	0.0	0.0
13	0	1	0	615	110	4	20	750	0.0	0.001	0.0	<b>0.82</b>	0.147	0.005
14	881	5	5	5	3	0	22	921	<b>0.957</b>	0.005	0.005	0.005	0.003	0.0
15	0	1	0	40	1250	156	1	1448	0.0	0.001	0.0	0.028	<b>0.863</b>	0.108
16	124	89	65	42	213	24	205	762	0.163	0.117	0.085	0.055	0.28	0.031
17	1	733	7	0	6	0	29	776	0.001	<b>0.945</b>	0.009	0.0	0.008	0.0
18	72	105	112	24	38	7	25	383	0.188	0.274	0.292	0.063	0.099	0.018
19	1	502	3	0	0	0	1	507	0.002	<b>0.99</b>	0.006	0.0	0.0	0.0
20	0	720	1	0	0	0	1	722	0.0	<b>0.997</b>	0.001	0.0	0.0	0.0
21	11	13	630	1	1	0	5	661	0.017	0.02	<b>0.953</b>	0.002	0.002	0.0
22	248	297	209	15	7	3	14	793	0.313	0.375	0.264	0.019	0.009	0.004
23	396	0	3	620	335	4	3	1361	0.291	0.0	0.002	0.456	0.246	0.003
24	1	681	1	0	0	0	5	688	0.001	<b>0.99</b>	0.001	0.0	0.0	0.0

482  
483

**Appendix 2** - The result of the clustering for ~26,000 frontal images with k=25. Values above the threshold of 0.7 are shown in green.