

# Creating an Additional Class Layer with Machine Learning to counter Overfitting in an Unbalanced Ancient Coin Dataset

Sebastian Gampe\*<sup>1</sup>, Karsten Tolle<sup>1</sup>

<sup>1</sup> Big Data Lab, Goethe-Universität – Frankfurt am Main, Germany

\*Corresponding author

Correspondence: gampe@em.uni-frankfurt.de

## ABSTRACT

We have implemented an approach based on *Convolutional Neural Networks* (CNN) for mint recognition for our *Corpus Nummorum* (CN) coin dataset as an alternative to coin type recognition, since we had too few instances for most of the types (classes). However, this shift increased an existing problem with our dataset: the extremely unbalanced number of instances per class. While some of our classes consist of only 20 instances, others consist of several hundred. After training our VGG16 model we unsurprisingly observed an overfitting of these “big” classes within the confusion matrix. To reduce this problem, we tried to split the dominating classes with the most images into several smaller ones and called them additional class layers. We use three different machine learning (ML) approaches to perform this breakdown. One is an unsupervised clustering method without additional manual work. The other two are supervised approaches which explicitly take into account the motifs of the coins themselves: a) an Object Detection model that predicts trained entities, and b) a *Natural Language Processing* (NLP) method to find entities in the textual descriptions of the coins. Based on the combination of obverse and reverse results from these two approaches new additional class layers were defined for each of them independently. After retraining our mint recognition model with these new classes, we evaluated the results based on the confusion matrix. In our case, the best results could be observed by forming an additional class layer based on the NLP method. Unfortunately, in our situation the overfitting problem could only be reduced and not eliminated.

**Keywords:** Machine Learning, Image Recognition, Convolutional Neural Networks, Unbalanced Dataset, Ancient Coins

## Introduction

42 In our current project D4N4 – Data quality for Numismatics based on Natural language processing and  
 43 Neural Networks (D4N4 n.d.) we want to implement a Machine Learning (ML)-based coin type recognition  
 44 model that covers as many coin types as possible from the “Corpus Nummorum” (CN) (Corpus Nummorum  
 45 n.d.) dataset. The goal is to use it to improve and verify the data quality of existing data, and also use it to  
 46 help the process of entering new coins. The CN dataset features about 19,600 different coin types and  
 47 more than 49,000 coins from four different ancient landscapes (Thrace, Moesia Inferior, Troad and Mysia).  
 48 This dataset contains coins from several different museums, institutions and collectors. The largest part of  
 49 our images comes from the Berlin-Brandenburgische Akademie der Wissenschaften, the Münzkabinett  
 50 Berlin and the Bibliothèque nationale de France, Département des Monnaies, médailles et antiques. We  
 51 published the images from these three institutions (due to existing copyrights) as a dataset for ML research  
 52 on Zenodo (Corpus Nummorum 2023). A coin in the database is generally represented by images of the  
 53 obverse and reverse of the original, or by a plaster cast. In rare cases, both representations are assigned to  
 54 a coin record. For our ML dataset, we merged the obverse and reverse images of a coin into a single image  
 55 showing both (as can be seen in fig. 1).

56 The biggest challenge for our type recognition is the ratio of an average of approximately two coin  
 57 images per type. In previous experiments we learned that 20 coin images per class threshold is a reasonable  
 58 starting point to achieve good results in the training for our data. This means that for most of the types our  
 59 dataset currently has too few coins (Gampe 2021, Gampe and Tolle in print 2019). Currently only 179 of  
 60 the 19,600 type classes can meet the condition. The VGG16 model, which was pretrained on the ImageNet  
 61 dataset (ImageNet 2021), achieves a Top-1 Accuracy of 82% (tab. 1) based on those types that meet the  
 62 threshold (>20 images), with the drawback that 99% of the type classes in the CN dataset are not part of  
 63 the training<sup>1</sup>. We are constantly increasing the number of our images and we have been able to double the  
 64 number of trainable types since the start of the project.



65

66 **Figure 1** – The merged images of CN coin 2377 (CN type 763; mint : Maroneia) and CN coin 18232  
 67 (CN type 11944; mint: Pergamon) with black bars for the quadratic input format of the CNNs (Photos :  
 68 Münzkabinett Berlin)

69 In order to get better coverage and still generate something useful, we trained a different model to  
 70 recognize another important aspect of our coins: the mint in which the coin was produced. Predicting a  
 71 coin’s mint could also reduce the workload for our numismatists since the presorting of a larger number of  
 72 coins by mint can save a lot of time. We kept our VGG16 setup and only needed to change the training and  
 73 test set. Currently there are 122 mints in our data set. With the 20-coin image threshold, 98 of them are  
 74 eligible for training. This way we could use about 40,000 (~80%) of our images for this approach. The  
 75 remaining 24 mints have less than 200 images combined. The accuracy values here are comparable to the  
 76 type recognition (Top-1 Accuracy: 79%, tab. 1) in spite of the fact that the individual mint classes are much  
 77 more disparate than the type classes. Most mint classes consist of several different coin types which differ

<sup>1</sup> The Top-1 Accuracy shows the percentage of correct results of the predictions with the highest probability. The Top-5 Accuracy is the percentage of correct predictions within the five most likely predictions.

78 more or less from each other, Pergamon for example has 653 different coin types (630 are currently  
 79 published on the website)<sup>2</sup>.

80 **Table 1** – Metrics for the models with and without the additional class layers

CNN Architecture	Class Type	Additional Class Layer	Number of Perinthos / Pergamon classes	Top-1 Accuracy	Top-5 Accuracy
VGG16	Types	None	1	82%	98%
VGG16	Mints	None	1	79%	94%
VGG16	Mints	DeepCluster	15 / 15	73%	91%
VGG16	Mints	DeepCluster	10 / 10	74%	92%
VGG16	Mints	Object Detection	8 / 10	78%	93%
VGG16	Mints	Object Detection	4 / 4	77%	93%
VGG16	Mints	NLP	16 / 16	76%	92%
VGG16	Mints	NLP	8 / 9	78%	93%

81  
 82 However, we still encounter the problem of having an unbalanced dataset when training the mint  
 83 recognition. The advantage of the significantly higher number of usable images comes with the problem of  
 84 a clearly higher amount of instances in a single class. Due to the different production patterns of the mints  
 85 and the disparity in the number of coins of each ancient city in our area of interest, we also have  
 86 significantly different coin image numbers for each city. Some of them have only the threshold value of 20  
 87 images while others have several hundred (fig. 2).

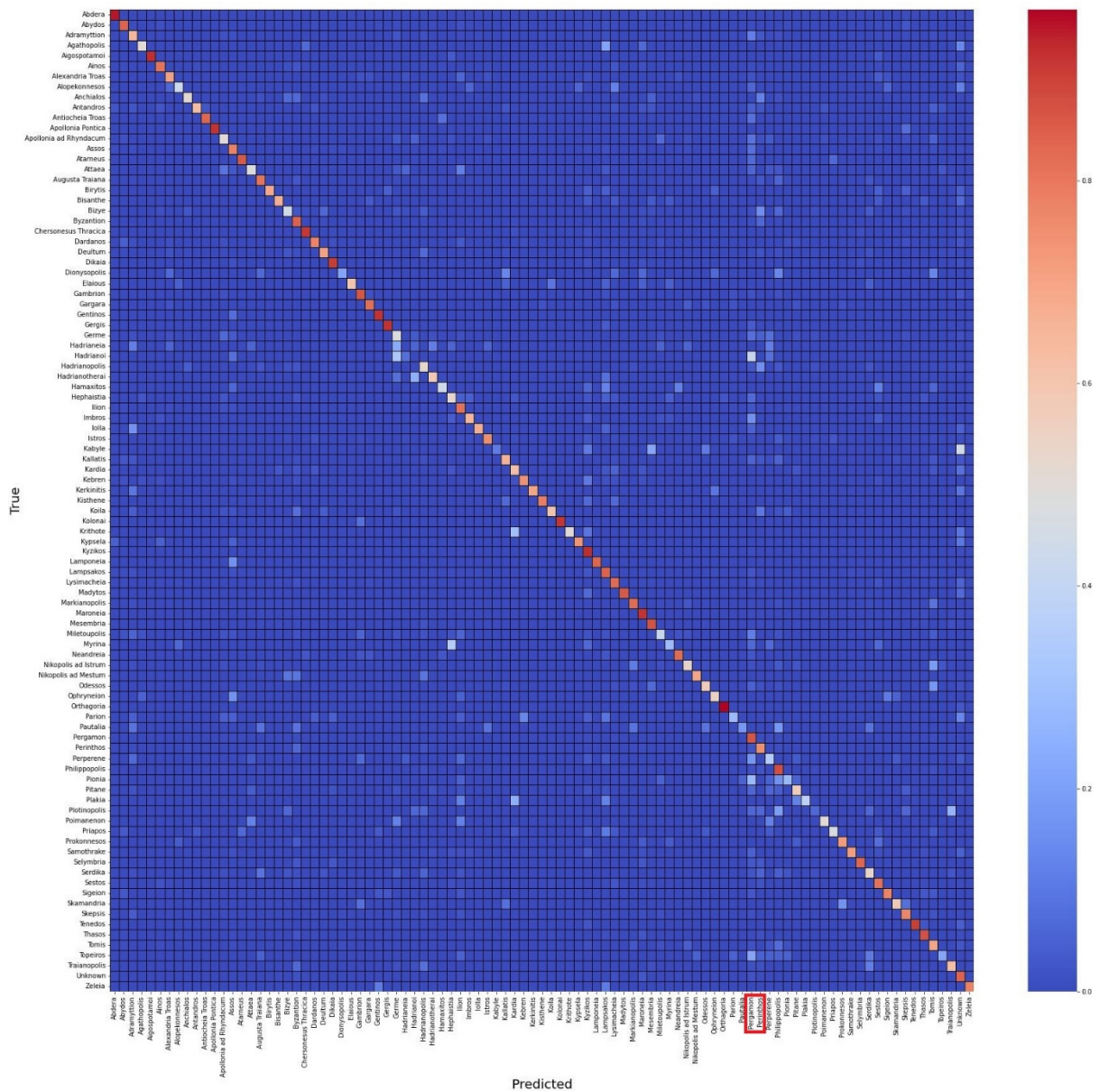
88 Such big differences in the number of images per class can lead to a phenomenon known as overfitting.  
 89 A CNN model learns to recognize some classes better than others due to the larger number of images, with  
 90 the result that the weights in the trained network can be more tuned to these classes during training.  
 91 Subsequently, it is possible that images of other smaller classes are predicted preferentially as belonging  
 92 to one of these overfitted classes. Overfitting can also appear by too many epochs during the training. In  
 93 our case and in this paper this source for overfitting is neglected and we concentrate on the inhomogeneity  
 94 of the data set as a basis for overfitting. A good way to determine if overfitting is present in a model is to  
 95 create the so-called Confusion Matrix (CM) (fig. 3). The CM is a common way to make overfitting visible for  
 96 individual classes in a multiclass problem. In the CM the predictions of the model and the true classes  
 97 (ground truth) are confronted. In the case of a model that is 100% correct for all predictions, the diagonal  
 98 of this matrix should be deep red. Recurring errors for one class create visible points off the diagonal.  
 99 Traces of vertical lines for one or more classes are a good indicator for the overfitting problem. The CM of  
 100 our most recent model shows these traces (fig. 3). They are clearly visible for the cities that have the most  
 101 associated images like Pergamon and Perinthos (fig. 2). Every city with more than 1500 assigned images  
 102 shows clear signs of overfitting on the CM. This overfitting is probably also responsible for the poorer  
 103 performance of the model. The Pergamon line is the most prominent. Although Perinthos also has a large  
 104 number of images, the overfitting problem is less pronounced here than with Pergamon. These differences  
 105 and the fact that both mints produced very different looking coin types and thus present a challenge when  
 106 dividing into new smaller classes make them a good case study in our view.

107 One idea for solving the overfitting problem is limiting the number of training images per class. The  
 108 problem for the image limitation is the number of types which belong to a mint. For example, for the mint  
 109 Pergamon there are over 3600 coin images, which are split over 653 different coin types. This means that  
 110 if we want to represent each type with at least one coin image in the training set the Pergamon class would  
 111 still be significantly larger than many other classes. A limited dataset can also lead to lower model

<sup>2</sup> Overview of the 630 published types in the CN dataset:  
<https://www.corpus-nummorum.eu/search/types?type=quicksearch&mints%5B%5D=74> [accessed 15 December 2023].



128 the improvement of machine learning (ML) algorithms as such, but on the setup, in particular the handling  
 129 of the datasets and the definition of the result classes for the training.



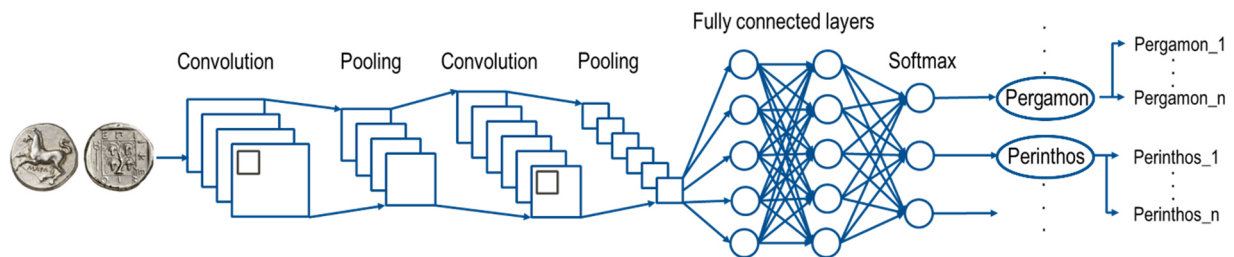
130  
 131 **Figure 3** – Confusion matrix for the mint model without the additional class layer. (Graphic: S. Gampe,  
 132 Big Data Lab)

133 The goal was to improve our mint recognition model by breaking down large classes with many input  
 134 images into several smaller ones. For this purpose we used three different methods: 1. DeepCluster – an  
 135 unsupervised clustering method, 2. Object Detection based on a Region Based – Convolutional Neural  
 136 Network (R-CNN) and 3. our Natural Language Processing (NLP) pipeline (Gampe and Tolle in print 2019).  
 137 The DeepCluster model is available on GitHub (GitHub – facebookresearch / deepcluster n.d.). The Object  
 138 Detection approach trained for this paper is based on TensorFlow and Keras (Keras n.d., TensorFlow n.d.).  
 139 The already existing NLP Pipeline was developed with the spaCy application programming interface (spaCy  
 140 n.d.). Our Image Recognition models, which are based on a pretrained VGG16 model, are implemented  
 141 with TensorFlow and Keras (Gampe and Tolle in print 2019, TensorFlow n.d., Keras n.d.). All of the above  
 142 methods run on Jupyter Notebook and are written in Python programming language.

143 In another Project (ClareNet) the DeepCluster Method was used to cluster a coin hoard with celtic  
 144 coins. The already existing typology and the allocation of coins to them was checked with the method. The

145 R-CNN based Object Detection was utilized to crop the area of a portrait on imperial roman coins to prevent  
146 a CNN from making decisions based on the legend (Gampe 2021). Our NLP approach has been in  
147 development for some time (Klinger et al. 2018, Gampe and Tolle in print 2019) and is now also used for  
148 other reasons within the CN-project. We therefore wanted to check whether it could also be used to solve  
149 our overfitting problems.

150 Our idea for addressing this problem is to break down big classes like Pergamon and Perinthos into  
151 smaller ones. We call this the *additional class layer*. The term “layer” is well known from neural networks  
152 like the CNN, which has different layers (fig. 4). The last layer in such networks is called the softmax and is  
153 responsible for transforming the incoming numerical values from the preceding layers into the class’s  
154 probabilities (Amidi and Amidi n.d.). We call this layer also the class layer. Our idea is to add an additional  
155 layer for large classes on top of the existing ones (fig. 4).



156

157 **Figure 4** – Overview of a convolutional neural Network, its different layers and the additional class  
158 layer. (Photo: Corpus Nummorum. Graphic: S. Gampe, Big Data Lab)

159 In this process, classes like Pergamon are divided based on the similarity of their different types and  
160 new smaller classes are added. Then we add the old class layer with the new ones to get a new train and  
161 test set. To realize this approach, we use three different Machine Learning based methods:

- 162 1. Unsupervised “DeepCluster”
- 163 2. Region Based – CNN based Object Detection
- 164 3. Natural Language Processing

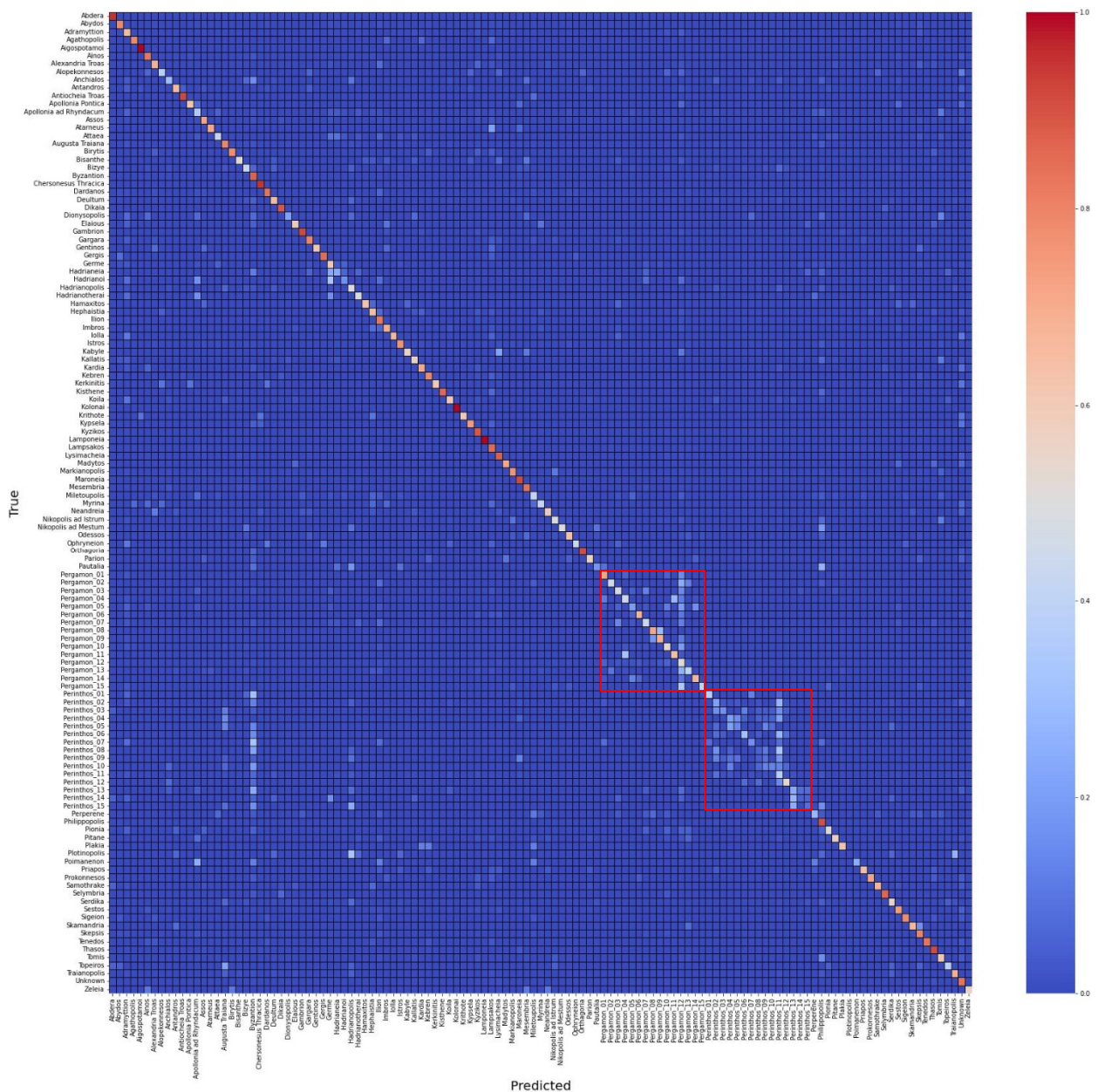
165 By creating these new classes and reducing the old large ones we tried to reduce the amount of overfitting  
166 and potentially increase the accuracy of our models. To test these approaches, we applied them to two  
167 different mints: Pergamon with c. 3600 images and Perinthos with c. 1800 images. Pergamon was chosen  
168 because it is the mint with the most associated images. In contrast Perinthos has just half as many images,  
169 however, this city has types that differ greatly from those from Pergamon due to the fact that the city  
170 belongs to another region (Pergamon – Mysia, Perinthos – Thrace). These differences are important to  
171 ensure the applicability of the three approaches to different mints.

172

### Unsupervised “DeepCluster”

173 Our first method, *DeepCluster* from Facebook Research (Caron et al. 2018), combines unsupervised and  
174 supervised elements. An integrated CNN extracts features from the input images. Afterwards, these  
175 features are clustered with a k-means algorithm and the resulting clusters are used as pseudo labels. These  
176 newly labeled images serve as input for further CNN training. The number of clusters depends on the input  
177 parameter “k”, which can be freely selected before the start (Caron et al. 2018). One big benefit of this  
178 approach is the fact that DeepCluster needs no time-consuming adaptations to our problem. Values for a  
179 few parameters have to be chosen, such as the number of clusters “k” and the amount of training epochs.  
180 We chose 200 epochs for each attempt. The number of clusters was set to 15 and 10 because we didn’t  
181 want to create a large number of small classes. DeepCluster’s Algorithm uses all of our images from  
182 Pergamon and from Perinthos to form the 15 clusters for each mint separately. This means that both  
183 original classes have been completely split up and therefore no longer exist as classes in the training.  
184 However, it is possible that DeepCluster forms very inhomogeneous clusters with a lot of different looking  
185 types.

186 For our first attempt of dividing the Pergamon and Perinthos classes into smaller ones we had 15  
 187 clusters each (Pergamon\_01 to Pergamon\_15 and Perinthos\_01 to Perinthos\_15). Although both classes  
 188 have different numbers of images, we have chosen an equal number of new classes to investigate the  
 189 effects on these different sized classes. After the training most clusters had a size of 50 to 100 images for  
 190 Perinthos and 100 to 200 for Pergamon. Both mints also had one cluster with many more images than the  
 191 others: DeepCluster tends to build what we call “garbage clusters”, which contain all images that could not  
 192 be assigned to the other clusters (Pergamon\_10 and 12, Perinthos\_02 and 11). They are somewhat  
 193 comparable to the original classes with the leftover images from our two other approaches (see below).  
 194 The 30 clusters built with DeepCluster are now incorporated as new classes in our train and test set  
 195 replacing the original Pergamon and Perinthos classes. After the training we could observe that the Top-1  
 196 and Top-5 Accuracy values were below those from the unmodified model (tab. 1). We also created the  
 197 confusion matrix for the new model on the test set (fig. 5).



198

199 **Figure 5** – Confusion matrix for the DeepCluster approach. (Graphic: S. Gampe, Big Data Lab)

200 What is immediately noticeable is that there is a clear vertical line in the area of the new classes. This  
 201 means that these classes are often confused with each other. A look at the composition of the clusters

202 indicates that this is due to the inhomogeneity of the assigned images within each cluster. While many  
203 clusters share images of the same or very similar types, the two biggest “garbage clusters” from both mints  
204 (Pergamon\_12 and Perinthos\_11) are clearly visible on the CM. These share a particularly large number of  
205 coin types with the other new classes. For Pergamon, the overfitting problem has diminished, but some of  
206 the brighter points from the unmodified class columns are now spread across the vertical lines of the new  
207 classes. For Perinthos we could also observe a slight decrease of the overfitting. However, new overfitting  
208 problems can also arise, as it is apparent in the Byzantion column at the level of the new Perinthos classes.  
209 We repeated this test with a smaller number of clusters. This time we executed DeepCluster with 10  
210 clusters as preset. However, the result barely changed. The confusion between the new classes is still there  
211 and the overfitting for both mints is nearly the same as with 15 clusters. We also found a reduction of the  
212 Top-1 Accuracy of about 5% in both tests. These results have led us to conclude that the DeepCluster  
213 method is not suitable for our problem.

214

## 215 **Region Based – CNN based Object Detection**

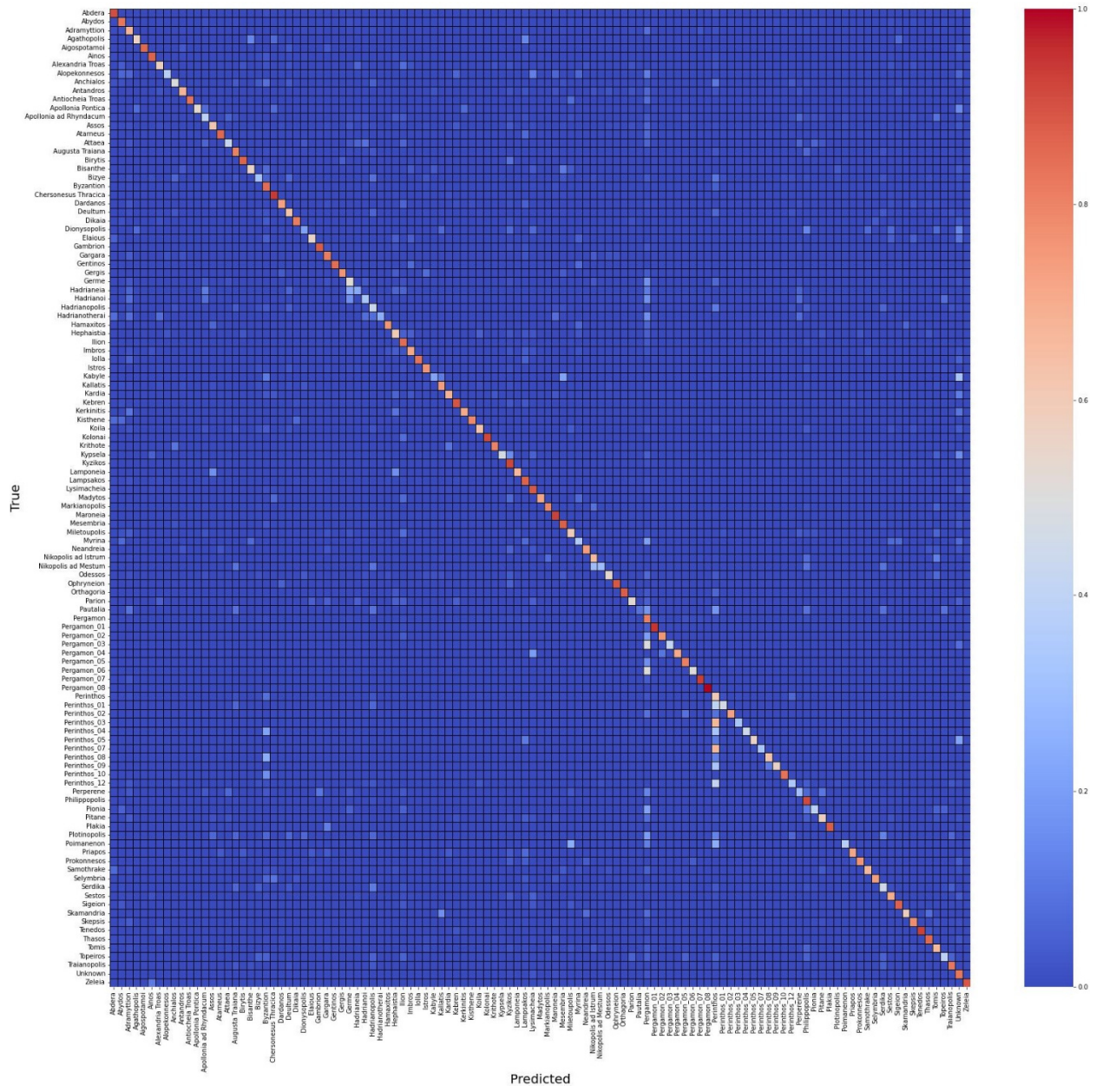
216 Our second approach is an Object Detection model from Keras (Keras n.d.). It is based on a Region Based  
217 Convolutional Neural Network (R-CNN) which produces a set of region proposals that are likely to contain  
218 objects, and uses a CNN to extract features from each region proposal to classify objects within these  
219 regions. (Girshick et al. 2014). The way new classes were created here followed a different concept. We  
220 trained the R-CNN model on frequently occurring subjects on the coins like “head” or “sitting person”.  
221 Most of these subjects are among the most common objects and animals in the CN database (Wirth 2021).  
222 The new classes were built based on the combinations of these subjects in the dataset<sup>3</sup>. The training of the  
223 R-CNN model was carried out by Huy Long, who wrote his master’s thesis on this topic (Long 2022). He  
224 annotated 20 to 30 images for every subject with a polygonal annotation and with bounding boxes. With  
225 a polygonal annotated training set the model can predict the contour of an object, but annotation is more  
226 time consuming compared to the annotation of bounding boxes. After training the Object Detection model,  
227 it was used to predict subjects on each coin image of Pergamon and Perinthos. Based on these results, new  
228 classes had been built manually based on the resulting combinations of obverse and reverse subjects (e.g.  
229 Head–Owl). For Pergamon this generated eight and for Perinthos ten classes. These classes have different  
230 sizes with around 20 to 200 images assigned to them. This way we reduced the number of images in the  
231 original Pergamon and Perinthos classes to 2,520 and 1,170. After training our VGG16 model on the  
232 updated train and test set we observed that the top-1 and top-5 accuracy has hardly changed (tab. 1). But  
233 the resulting CM shows a problem similar to the DeepCluster approach (fig. 6). Images of the new classes  
234 are often attributed to the reduced Pergamon and Perinthos class. Furthermore, the overfitting for both  
235 original classes was not reduced as the CM shows. We repeated this test with a reduced number of new  
236 classes, where the images of some of the smaller classes had been manually merged back to the Pergamon  
237 and Perinthos classes. After retraining, beside a small decrease in the accuracy values, the overfitting  
238 problem also remained.

239 Clearly this approach is unfortunately not suitable for solving, or at least reducing our problem. An  
240 explanation for this could be the performance of the R-CNN Object Detection model. When examining the  
241 new classes, it became apparent that quite a number of objects had not been detected by the Object  
242 Detection model. The new classes were not really well distinguishable from the Pergamon and Perinthos  
243 ones. The annotation used seems not to cover the whole range of several subjects.

---

<sup>3</sup> The following subjects were used for the combinations: 1. Pergamon: “sitting\_person”, “head”, “serpent\_box”, “owl”, “serpents”, “eagle”, “podium”, “bull”, “bow”. 2. Perinthos: “sitting\_person”, “head”, “podium”, “bull”, “quadriga”, “double\_horse”, “club”, “laurel\_wreath”, “price\_crown”, “ship”, “table”, “pot”, “standing\_person”.





244  
245 **Figure 6** – Confusion matrix for the Object Detection approach. (Graphic: S. Gampe, Big Data Lab)

246  
247 **Natural language processing**

248 Our final method for dealing with mint overfitting is the application of our Natural Language Processing  
 249 (NLP) pipeline. We already have trained an NLP model for the textual descriptions of our coins in the CN  
 250 database (Gampe and Tolle in print 2019, Klinger et al. 2018). This model is able to find numerous entities  
 251 like “Athena” or “Spear”. The list of trained entities includes four different categories: persons, objects,  
 252 animals and plants. We wrote a query for grouping all coins with the same entities on the obverse and  
 253 reverse, and sorted them by the frequency of these individual combinations. To create the new classes, we  
 254 had to filter these combinations manually because a coin’s description can have several entities assigned  
 255 to it and overlaps are possible. The example (fig. 7) shows all entities in a coin’s obverse and reverse  
 256 description. For example, this coin can be assigned to the combination “Athena–Owl” and “Head–Palm

257 Branch". Filtering is a time-consuming but mandatory step due to the large number of combinations. We  
258 also wanted to avoid new classes sharing coin images.



259

260 **Figure 7** – Coin image descriptions and entities found by the NLP model. These descriptions are used  
261 for different coins and/or coin types (fig. 1), e.g. coins [https://www.corpus-](https://www.corpus-nummorum.eu/coins/53267?lg=en)  
262 [nummorum.eu/coins/53267?lg=en](https://www.corpus-nummorum.eu/coins/53267?lg=en) and [https://www.corpus-](https://www.corpus-nummorum.eu/coins/53408?lg=en)  
263 [nummorum.eu/coins/53408?lg=en](https://www.corpus-nummorum.eu/coins/53408?lg=en)  
264 [accessed 15 December 2023] for the Athena description shown here. (Graphic: S. Gampe, Big Data Lab)

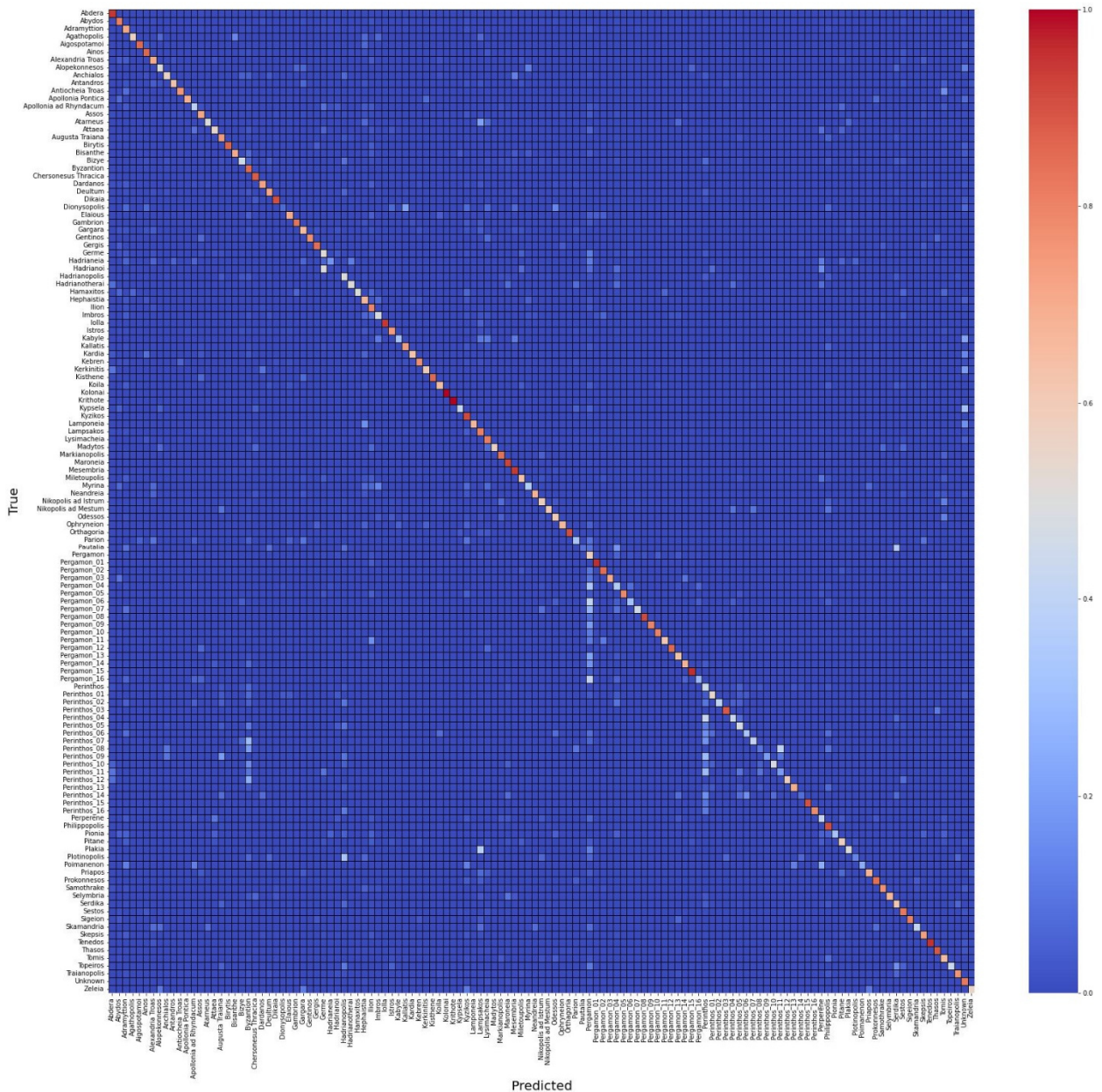
265 After the filter step the new classes could be built from the remaining combinations. We assembled 16  
266 new classes for Pergamon and Perinthos with different combinations of entities<sup>4</sup>. These classes contained  
267 mostly different types that share a similar appearance. The number of images for each class varied from  
268 40 to 400 for Pergamon, and 40 to 300 for Perinthos. Using the NLP model, we were able to reduce the  
269 number of images in the original Pergamon and Perinthos classes to 1201 and 629. This is a significantly  
270 larger reduction than that with Object Detection. Compared to the original model, the values of the metrics  
271 after the training are only slightly lower (tab. 1).

272 However, as the confusion matrix shows, the new classes are again often confused with the original  
273 Pergamon and Perinthos classes (fig. 8). The overfitting of Pergamon was slightly reduced with this attempt,  
274 but the vertical line in the original Pergamon class fields (outside the new classes) is still visible. The  
275 reduction of the overfitting for Perinthos, on the other hand, is clearly visible. Due to the confusion of the  
276 old with the new classes both overfitting reductions had no influence on the metrics. We repeated this  
277 experiment again with a lower number of new classes (eight for Pergamon and nine for Perinthos).  
278 However, this had no positive effect on the overfitting problem for Pergamon. In fact, it became worse  
279 than before. The overfitting of Perinthos remained at the same level as in the first NLP attempt.

280 For this approach we can say that the NLP method required the most manual work of all three  
281 approaches. Our efficient NLP pipeline helped us best to separate the types that share a similar appearance  
282 for the new classes from the original classes. It also gave the best results for the overfitting problem based  
283 on the observation of the CM. However, the confusion between the new classes shows that the types of  
284 one mint very likely share some common features. This is something that could be further explored.

---

<sup>4</sup> The following entities were used for the combinations: 1. Pergamon: "Athena", "laurel\_wreath", "owl", "Augustus", "crepidoma", "bust", "Asclepius", "emperor", "Telesphorus", "cista", "serpent", "figure", "head", "bow", "heads", "temple", "trophy", "paludamentum", "Zeus". 2. Perinthos: "bust", "ears\_of\_corn", "Herakles", "lyre", "palm\_branch", "patera", "torch", "cuirass", "apples", "head", "altar", "athlete", "club", "Dionysos", "horses", "patera", "Isis", "apis", "radiate\_crown" "board".



285

286

**Figure 8** - - Confusion matrix for the NLP approach. (Graphic: S. Gampe, Big Data Lab)

287

288

### Summary and Conclusion

289

Due to our limited dataset and the high number of types, a training on types with sometimes only two or three instances per type results in an unacceptable performance (34% Top-1 Accuracy). By switching to mints (accumulating all types of a mint to one training class), a very unbalanced dataset was generated with some mints (classes) dominating the training. While the performance of our trained models seems sufficient, we saw traces of overfitting for the dominating mints in the confusion matrix. The problem of unbalanced data sets is very common in the domain of numismatics and archaeology in general. The overall preservation quality of the coins we used was relatively good, especially compared to excavation finds where additional problems and uncertainties occur. However, our work could lay the basic foundation to expand the models also to less well preserved material. It is therefore important to solve problems before dealing with less optimal material.

298

299 We conducted three different machine learning based experiments to solve this overfitting problem  
300 with the unbalanced Corpus Nummorum dataset. To do this, we split two mint classes with a large number  
301 of images and created an additional class layer for them, and generated another training and test set with  
302 it. In our first attempt we created several new classes based on DeepCluster (unsupervised). The generated  
303 clusters contained too many different looking coin types (based on a human judgment) which negatively  
304 affected the learning process of our VGG16 model. In the second attempt the class layer was created with  
305 an Object Detection approach. This generated only a few smaller extra classes and the remaining coins  
306 without objects detected still formed a dominating class. It must be stressed that the Object Detection  
307 approach was only trained for some very common objects and the overall performance was still limited.  
308 Both approaches did not produce appropriate solutions for our problem. In the third approach classes were  
309 generated with Natural Language Processing. They were the most distinguishable from the original classes  
310 and this approach reduced overfitting the most. However, due the amount of manual work and the  
311 confusion of new and old classes, we are currently not following this path either. The accuracy of all newly  
312 trained models was below the original mint model. Even the observed visual reduction of the overfitting  
313 for the original Pergamon and Perinthos classes in the NLP experiment had no positive impact on model  
314 performance. This means that so far we could not solve our overfitting issue with mint prediction in a  
315 sufficient way.

316 These approaches might be useful methods in other cases, however, it shows that a generic approach  
317 to an overfitting due to dominant classes has not been found. We are therefore investigating additional  
318 ways to tackle this problem. This includes the creation of new coin images using ML-based methods for the  
319 individual classes. It can be seen as an augmentation approach for the smaller classes in order to reduce  
320 the domination of some huge classes. Furthermore, we started to compare other model approaches than  
321 just CNN (VGG16), like vision transformers or multimodal approaches.

322 Our next steps in the D4N4 project are:

323 The domain experts are trying to improve the CN dataset by including more images, especially for smaller  
324 classes. We are also working on a Generative Adversarial Network (GAN) approach to create virtual new  
325 coin images (that never existed) for those classes with very few coins. Finally, we published our CN dataset  
326 for other scientists and students to test their own ML methods, for example those that were recently part  
327 of a data challenge course at the Goethe-University: (Corpus Nummorum 2023).

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

329 Google Colab notebook for testing our type and mint model is available online:  
330 <https://github.com/Frankfurt-BigDataLab/IR-on-coin-datasets>

331 NLP pipeline code, models and a Google Colab Notebook for testing are available online:  
332 <https://github.com/Frankfurt-BigDataLab/NLP-on-multilingual-coin-datasets>

## 333 **Conflict of interest disclosure**

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

## 336 **Funding**

337 The D4N<sup>4</sup> project is funded by the Deutsche Forschungs Gemeinschaft (DFG) in the program “e-  
338 research-Technologien”.

## 339 **References**

340 Amidi, Afshine, Shervine Amidi. *Convolutional Neural Networks cheatsheet*. n.d.  
341 <<https://stanford.edu/~shervine/teaching/cs-230/cheatsheet-convolutional-neural-networks>>  
342 [accessed 8 August 2023]

343 Caron, Mathilde, Piotr Bojanowski, Armand Joulin, Matthijs Douze. 2018. *Deep Clustering for Unsupervised*  
344 *Learning of Visual Features*. <<https://doi.org/10.48550/arXiv.1807.05520>> [accessed 8 August 2023]  
345 Corpus Nummorum. n.d. *Corpus Nummorum* <<https://www.corpus-nummorum.eu>> [accessed 8 August  
346 2023]  
347 Corpus Nummorum 2023. *Corpus Nummorum - Coin Image Dataset*.  
348 <<https://zenodo.org/doi/10.5281/zenodo.10033992>> [accessed 17 November 2023]  
349 D4N<sup>4</sup>. n.d. *Data quality for Numismatics based on Natural language processing and Neural Networks*.  
350 <<http://www.bigdata.uni-frankfurt.de/d4n4>> [accessed 8 August 2023]  
351 focal-loss. n.d. *focal\_loss.sparse\_categorical\_focal\_loss*.  
352 <<https://focal->  
353 [loss.readthedocs.io/en/latest/generated/focal\\_loss.sparse\\_categorical\\_focal\\_loss.html](https://focal-loss.readthedocs.io/en/latest/generated/focal_loss.sparse_categorical_focal_loss.html)> [accessed 8  
354 August 2023]  
355 Gampe, Sebastian. 2021. *Neuronale Netze zur Bestimmung römischer Kaiser auf Bildern antiker Münzen*  
356 (master thesis, Goethe-Universität Frankfurt a. M.). <[http://www.bigdata.uni-frankfurt.de/wp-](http://www.bigdata.uni-frankfurt.de/wp-content/uploads/2022/05/Masterarbeit_Sebastian_Gampe_online.pdf)  
357 [content/uploads/2022/05/Masterarbeit\\_Sebastian\\_Gampe\\_online.pdf](http://www.bigdata.uni-frankfurt.de/wp-content/uploads/2022/05/Masterarbeit_Sebastian_Gampe_online.pdf)> [accessed 8 August 2023]  
358 Gampe, Sebastian, Karsten Tolle. in print 2019. *Combination of Machine Learning Methods of Image and*  
359 *Natural Language Recognition of Ancient Coin Data, Computer Applications & Quantitative Methods in*  
360 *Archeology (CAA), Proceedings of the conference in Krakow 2019*  
361 GitHub - facebookresearch / deepcluster. n.d. *Deep Clustering for Unsupervised Learning of Visual Features*.  
362 <<https://github.com/facebookresearch/deepcluster>> [accessed 8 August 2023]  
363 ImageNet. 2021. *ImageNet*. <<https://www.image-net.org/update-mar-11-2021.php>> [accessed 8 August  
364 2023]  
365 Keras. n.d. *Keras: Deep Learning for humans*. <<https://keras.io>> [accessed 8 August 2023]  
366 Klinger, Patricia, Sebastian Gampe, Karsten Tolle, Ulrike Peter. 2018. *Semantic Search based on Natural*  
367 *Language Processing: a Numismatic example. Journal of Ancient History and Archaeology (JAHA)*, 5.3  
368 68–79. <<https://doi.org/10.14795/j.v5i3.334>> [accessed 8 August 2023]  
369 Long, Hui. 2023. *Klassifizierung von Motiven auf antiken Münzen mit Mask R-CNN* (master thesis, Goethe-  
370 Universität Frankfurt a. M.). <[http://www.bigdata.uni-frankfurt.de/wp-](http://www.bigdata.uni-frankfurt.de/wp-content/uploads/2023/02/Masterthesis-Huy-Luong.pdf)  
371 [content/uploads/2023/02/Masterthesis-Huy-Luong.pdf](http://www.bigdata.uni-frankfurt.de/wp-content/uploads/2023/02/Masterthesis-Huy-Luong.pdf)> [accessed 8 August 2023]  
372 Selvaraju, Ramprasaath R., Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, Dhruv  
373 Batra. 2019. *Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization*.  
374 <<https://doi.org/10.1007/s11263-019-01228-7>> [accessed 8 August 2023]  
375 Girshick, Ross, Jeff Donahue, Trevor Darrell, Jitendra Malik. 2013. *Rich feature hierarchies for accurate*  
376 *object detection and semantic segmentation*. <<https://doi.org/10.48550/arXiv.1311.2524>> [accessed 8  
377 August 2023]  
378 scikit-learn. n.d. *sklearn.utils.class\_weight.compute\_class\_weight*.  
379 <<https://scikit->  
380 [learn.org/stable/modules/generated/sklearn.utils.class\\_weight.compute\\_class\\_weight.html](https://scikit-learn.org/stable/modules/generated/sklearn.utils.class_weight.compute_class_weight.html)>  
381 [accessed 8 August 2023]  
382 spaCy. n.d. *spaCy: Industrial-strength Natural Language Processing in Python*. <<https://spacy.io>> [accessed  
383 8 August 2023]  
384 TensorFlow. n.d. *Tensorflow Object Detection API*.  
385 <[https://github.com/tensorflow/models/tree/master/research/object\\_detection#tensorflow-object-](https://github.com/tensorflow/models/tree/master/research/object_detection#tensorflow-object-detection-api)  
386 [detection-api](https://github.com/tensorflow/models/tree/master/research/object_detection#tensorflow-object-detection-api)> [accessed 8 August 2023]  
387 Wirth, Alicia. 2021. *Einfluss von Bildannotationstechniken auf die Genauigkeit von Machine Learning-*  
388 *Modellen zur Objekterkennung* (research project, Goethe-Universität Frankfurt a. M.)  
389 <[http://www.bigdata.uni-frankfurt.de/wp-content/uploads/2021/12/2021\\_FP\\_Alicia\\_online.pdf](http://www.bigdata.uni-frankfurt.de/wp-content/uploads/2021/12/2021_FP_Alicia_online.pdf)>  
390 [accessed 18 November 2023]