
Reproducing Visual Explanations of Variational Autoencoders

Anonymous Author(s)

Affiliation

Address

email

Reproducibility Summary

1

2 **Scope of Reproducibility**

3 In this work we perform a replication study of the paper “Towards Visually Explaining Variational Autoencoders”.
4 This paper claims to have found a method to provide visual explanations of Variational Autoencoders (VAEs). The
5 paper’s primary claim is that their proposed method can generate gradient-based attention maps from the latent space of
6 a VAEs. This is visually demonstrated on the MNIST dataset. Moreover, these attention maps are claimed to be useful
7 for anomaly detection, which is demonstrated on the UCSD Ped1 and MVTec-AD datasets. Finally, this method is
8 integrated into a loss function to obtain the attention disentanglement loss. This loss is shown to improve latent space
9 disentanglement when integrated into a FactorVAE model, which is demonstrated on the dSprites dataset. This paper
10 aims to reproduce all of the claims stated above.

11 **Methodology**

12 In order to produce attention maps which can localize anomalies for the MNIST dataset the code from the repository of
13 the authors could be reused. Additional models were implemented for the UCSD Ped1 and MVTec-AD datasets based
14 on supplementary material provided for the original paper. To reproduce the latent space disentanglement results, a
15 AD-FactorVAE model – which combines the attention disentanglement loss and FactorVAE model – was implemented
16 based on the original paper, the original paper’s supplement, the FactorVAE paper and external code sources.

17 **Results**

18 The attention maps for the MNNIST dataset were successfully replicated. Nevertheless, we failed to reproduce the
19 results for the UCSD Ped1 and MVTec-AD dataset. Furthermore, we were unable to reproduce the results for the
20 AD-FactorVAE. Potential explanations for this could be the incorrect aggregation and/or weighting of the attention
21 disentanglement loss, not training the models for enough epochs, or not using the correct method to produce the
22 attention maps.

23 **What was easy**

24 Reproducing the attention maps and anomaly localization results for the MNIST dataset was relatively easy.

25 **What was difficult**

26 Reproducing the anomaly localization results for the UCSD Ped1, and MVTec-AD datasets was difficult. Additionally,
27 reproducing the disentanglement results using the AD-FactorVAE proved to be problematic.

28 **Communication with original authors**

29 Except for a reply to our initial e-mail updating them of our attempt to reproducing the paper and asking if they would
30 share more of their code, which they declined, the authors of the paper did not respond to any of our questions.

31 1 Introduction

32 This paper describes our attempt to reproduce the paper “Towards Visually Explaining Variational Autoencoders” by Liu
33 et al [1]. Their paper introduces a method to provide visual explanations for variational autoencoders (VAEs) [2], which
34 is inspired by the recent development of gradient-based attention maps for Convolutional Neural Networks (CNNs) that
35 aid in visualizing and understanding these models [3]. The attention maps highlight areas of an image that are important
36 for the classification (in the CNN case) or reconstruction (in the VAE case) of the respective image. It is important
37 to search for these explanations since deep learning models tend to be difficult to interpret due to their complicated
38 structure [4]. This complexity allows for superior performance but also prevents creators from understanding the reasons
39 behind their models’ outputs. If the models’ underlying assumptions cannot be made explicit, there is a possibility
40 that the predictions of the models rely on false assumptions. This could cause failure in later stages of deployment.
41 Furthermore, adequate explanations for predictions can improve consumers’ trust in the abilities and fairness of a model
42 [5]. Additionally, we could learn new properties of the data with these explanations, and improve our own decision
43 making abilities [6]. Deep models have been especially successful in the field of computer vision and have been widely
44 adopted in real-life tasks which makes the demand for proper understanding of their performance all the more urgent
45 [7, 8, 9]. The performance of techniques to produce explanations can generally only be determined by qualitative
46 analysis, which is highly subjective and selective since not all train or test instances can be evaluated. Thus, we decided
47 to reproduce a paper that not only promises to provide a method that can generate visual explanations for VAEs but also
48 states that these attention maps produce state-of-the-art anomaly localization results and can be used to improve latent
49 space disentanglement, which can both be quantitatively measured.

50 2 Scope of reproducibility

51 The original paper proposes a gradient-based technique to produce visual attention maps for VAEs [1]. It claims
52 that these attention maps are not only useful as explanations of VAE predictions but can also be used to perform
53 anomaly localization, which is demonstrated on the MNIST, UCSD Ped1 and MVTEC-AD datasets, and latent space
54 disentanglement, which is demonstrated on the dSprites dataset. In this paper we try to reproduce the attention maps and
55 the results for these two practical use-cases. In order to determine whether the paper can be successfully reproduced,
56 the following claims were constructed that capture the most important statements and which can be either confirmed by
57 our research (in which case the reproduction was successful) or rejected (in which case our reproduction has failed).

- 58 • The method introduced in [1] is able to generate gradient-based attention maps from the latent space learned
59 by a VAE model. These attention maps denote characteristics of a set of images that are important to the model
60 and intuitively make sense, which is verified by qualitative analysis on the MNIST dataset.
- 61 • The attention maps are useful for anomaly localization, which is demonstrated on the MNIST data set, for
62 which only qualitative analysis is provided. Additionally, it is shown on the UCSD Ped1 dataset where it
63 obtains a better AUROC score than the *Vanilla-VAE* method that simply computes the difference between the
64 original and reconstructed image. Finally, the performance for anomaly localization is state-of-the-art, which
65 is shown on the MVTEC-AD dataset where it obtains AUROC and IOU scores that are better or similar to those
66 obtained by the models compared in the MVTEC-AD paper.
- 67 • The attention maps can be used to improve latent space disentanglement by combining them into a loss
68 function that can be added to the VAE loss during training. This is proven by incorporating the loss with
69 the FactorVAE which results in a higher disentanglement score for the dSprites data set than the baseline
70 FactorVAE model, while keeping a similar reconstruction loss.

71 3 Methodology

72 The following section discusses the implementation of the necessary functions to reproduce the claims stated in the
73 previous section.

74 3.1 Attention maps

75 The first claim of the original paper is that the proposed gradient-based attention maps can explain the inner workings
76 of a VAE model. Attention maps can be generated for each latent dimension separately. First, the image is forwarded
77 through the model to obtain its latent representation. Subsequently, the activation of each of the dimensions of this
78 representation is backpropagated to one of the last convolutional layers. The resulting gradients are aggregated per
79 channel of the feature maps of the convolutional layer to obtain the weight given to each channel of the feature map.

80 This weighted sum of channels of the respective feature maps can then be used to retrieve the attention map of a single
81 latent dimension. The method to obtain attention maps for anomaly detection is also closely related to this technique.
82 The authors showcase how these attention maps explain the model by training models on a single class of images from
83 the MNIST data set [10] and then showing attention maps generated by the model for images of other classes. These
84 attention maps highlight areas which are not characteristic of the originally learned class.

85 The training of the VAE model was done by following the instructions in the repository with the code made available by
86 the authors which included the main training procedure and the experimental VAE architecture. Small adjustments
87 were implemented to make the core run on CPU as well. Acquiring the dataset, and transforming this to the right data
88 structure was already done in the provided code as well. In addition to training the models ourselves, we had access to
89 the pre-trained models which could be downloaded via a link in the authors' GitHub repository. The VAE structure –
90 two convolutional layers followed by a fully connected layer in the encoder and a mirrored structure in the decoder –
91 and the hyperparameters – `image_size = 28 | batch_size = 128 | latent_space = 32 | learning_rate 0.001` – were described
92 in the supplementary paper, which is available online¹, as well as provided in the code base. The amount of epochs was
93 unclear, but we trained the model for 100 epochs since this was the default provided in the code.

94 The provided code already contained an implementation for attention map generation, which we could run after
95 removing an erroneous parameter (`gradient=one_hot`) in the backpropagation function². However, when consulting the
96 paper it seemed that some parts of the code did not match the original method. By following the paper more closely we
97 discovered that Liu et al. distinguish between two different methods, (1) the first is proposed to visualise explanations
98 for VAEs by calculating the attention maps per latent dimension and aggregating these, and (2) the second is proposed
99 specifically for anomaly detection where attention maps are generated from the sum of the inferred mean vector. In the
100 provided code base only the (2) second method is implemented and so, in addition, we implemented the (1) first method
101 ourselves.

102 Besides implementing both methods, we also chose to test various functions in place of the ReLU in function (2) of
103 the paper. This was done as the provided code implemented the absolute value instead of the ReLU, which seemed to
104 give similar results and we were curious to see whether the sigmoid could be used as well, since this is also a common
105 function used to scale values between 0 and 1. To conclude, we tested both the pre-trained models and the models
106 trained from scratch with the two different methods and the three different functions.

107 **3.2 Anomaly detection**

108 The second claim by Liu et al. is that the attention maps are not only useful for anomaly localization but can produce
109 state-of-the-art results. They show this by conducting qualitative and quantitative experiments on two different datasets
110 for anomaly detection: the UCSD Ped1 dataset [11] and the MVTec-AD dataset [12].

111 **3.2.1 UCSD Ped1 set up**

112 The UCSD Ped1 dataset contains video frames of a pedestrian walkway where all non-pedestrians are considered to be
113 anomalies. The implementation for the anomaly detection on the UCSD Ped1 dataset was not provided in the authors'
114 code. For the qualitative research we could use the same approach as in the previous section, namely training a VAE
115 model on the UCSD Ped1 training set – which does not contain any anomalies – and afterwards generate attention maps
116 of the test set – which does contain anomalies – for each of the convolutional layers of the VAE using the (2) second
117 method. The changes we had to make for these experiments consisted of implementing a different architecture for the
118 VAE model containing three convolutional layers instead of two and adding functions to properly read in the UCSD
119 Ped1 training and test set, which included resizing the images to 100 by 100 pixels. Table 2 in the addendum paper
120 shows the architecture of the VAE model that was used for these experiments and specifies a learning rate of 0.0001
121 with a batch size of 32 frames for training. The amount of epochs was unclear, but we ended up training the model for
122 1000 epochs, which took around 4 hours on one NVIDIA 1080 GPU.

123 For replicating the quantitative results, we implemented functions to evaluate the produced attention maps. The UCSD
124 Ped1 test set contains both video frames containing anomalies as well as ground truth masks denoting the location of a
125 possible anomaly. Given the anomaly attention maps, binary anomaly localization maps were generated using a variety
126 of thresholds whose respective overlap with these ground truth masks is encapsulated in a pixel-level true positive rate
127 (TPR), false positive rate (FPR) and ROC curve. Subsequently, these could be used to calculate the area under the ROC
128 curve (ROC AUC) score using the sklearn metric '`roc_auc_score`'. Furthermore, with the trained model we tried to
129 reproduce the baseline score presented in the paper, which involved computing the difference between the input image
130 and its reconstruction and similarly comparing the result with the ground truth masks to obtain the ROC AUC score.

¹https://openaccess.thecvf.com/content_CVPR_2020/supplemental/Liu_Towards_Visually_Explaining_CVPR_2020_supplemental.pdf

²The original authors have recently also updated their implementation in the same way we did.

131 3.2.2 MVTec-AD set up

132 In the paper, the anomaly detection is further tested on the in 2019 released MVTec Anomaly Detection (MVTec-AD)
133 dataset [12]. This dataset consists of images of 15 natural objects and textures with different defects and pixel-level
134 ground truth masks. Again the repository did not contain any code to run experiments specifically for this dataset, but
135 the functions to obtain the attention maps using method (2) could be used. Additionally, we implemented the VAE
136 as stated in Table 3 of the addendum paper by the authors. The encoder consists of a ResNet18 (except for its last
137 two layers) followed by two linear layers and the decoder consists of two linear layers, six blocks containing a 2D
138 Convolutional layer, batch normalization and a ReLU, and finally a sigmoid layer. During training, Adam was used
139 for optimization with a learning rate of $1e-4$ and batch size of 8. We trained a separate model for each object/texture
140 category. Again, the number of epochs was undeclared, but we trained each category for 100 epochs due to time
141 constraints.

142 Before training, the data was augmented similarly to the MVTec-AD paper, of which the specifications are as follows.
143 All images were resized to 256×256 pixels. Additionally, they were randomly rotated between $[-30^\circ, +30^\circ]$ and/or
144 horizontally mirrored to create an augmented training set of 10.000 images per object/texture. The probability for the
145 horizontal mirroring was not specified but chosen to be 0.5 because it is the default value for this method. The padding
146 used after rotation was also unclear. We therefore decided to use black pixels since that is the standard parameter for
147 padding. For comparison, the *leather* category was padded differently, namely with the colour of the original image’s
148 upper-left pixel.

149 After training, the attention maps for the images containing anomalies could be generated for the last convolutional
150 layer of the VAE model. These attention maps could then be compared to the ground truth masks to obtain the ROC
151 AUC score. Moreover, the intersection-over-union (IOU) score was also calculated from the ROC curve.

152 The experiments were run on the same GPU node as the USCD data set. Creation of the augmented dataset took around
153 1,5 hours per object, while training of one model for 100 epochs took around 2,5 hours. The MVTec-AD dataset has a
154 size of 5GB and training requires around 11GB of memory

155 3.3 Latent space disentanglement

156 The authors’ code did not contain any implementation specific to the proposed attention disentanglement (AD) loss.
157 Nonetheless, the code to generate attention maps applying the (2) method proved to be useful. Consequently, we had
158 to implement the FactorVAE, disentanglement metric, and calculation of the AD loss ourselves. The authors use the
159 FactorVAE [13] to showcase how the proposed loss can be integrated into other loss terms. They call this new version
160 of the model AD-FactorVAE. To implement FactorVAE, we adapted an external GitHub repository [14] which was
161 recommended to us by the original authors. From this repository we also adopted code to read in the dSprites dataset,
162 which is used for the evaluation of both the FactorVAE and AD-FactorVAE [15]. We integrated the proposed loss
163 into this implementation according to Equation 6 in the paper. Besides the model, the paper also makes use of the
164 disentanglement metric proposed in the FactorVAE paper. The underlying assumption of this metric is that if the latent
165 space is perfectly disentangled, then each latent dimension of the representation should correspond to a single latent
166 factor used to generate the data. Consequently, if that latent factor is kept fixed, the variance of the corresponding
167 latent dimensions should be 0. A majority vote classifier is trained to predict factors corresponding to specific latent
168 dimensions by fixing single factors, passing data through the trained VAE model and checking for the dimension with
169 the lowest variance. The metric is the accuracy of this classifier. The implementation of this metric was adopted from
170 an external GitHub repository [16].

171 The hyperparameters used for these experiments are shown in Table 1. The original paper states that they did not modify
172 any parameters of the model compared the ones mentioned in the FactorVAE paper [13], other then increasing the
173 number of latent dimensions from 10 to 32. The negative slope of the LeakyReLU in the Discriminator is not mentioned
174 in the FactorVAE paper, however the adapted implementation assumes it to be 0.2. Moreover the authors also did not
175 specify the size of the test set used to obtain the disentanglement score. We chose to use 200 votes, because this did not
176 slow down the score calculation too much while still providing us with a robust indication of the accuracy. Additionally,
177 some hyperparameters of AD-FactorVAE were unclear from the paper. There is no optimal value stated for λ , the
178 weight of the AD loss. It was also not clear how to aggregate AD losses of several latent dimension pairs or how to
179 pick one or more pairs for the AD loss in general. However due to time constraints, we could only experiment with
180 the different aggregation of AD losses for all possible dimension pairings. The experiments were again run on one
181 NVIDIA 1080 GPU node. Training requires around 16GB memory and 24 hours to finish, while testing requires the
182 same amount of memory and at most 30 minutes.

Optimizer:	Adam	Learning rate:	10^{-4}	Batch size:	64
Latent dimensions:	32	VAE beta1:	0.9	VAE beta2:	0.999
LeakyReLU slope:	0.2	Discriminator beta1:	0.5	Discriminator beta2:	0.9
Metric batch size:	100	Metric training size:	800	Metric evaluation size:	200
γ value:	{20, 40, 100}	λ value:	1	Aggregation method:	{sum, mean}

Table 1: Hyperparameters used for AD-FactorVAE. **Bold** parameters are not clearly defined in the paper.

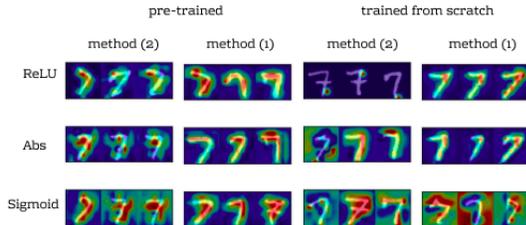


Figure 1: Qualitative results on the MNIST dataset with activation functions along the y axis and the various methods along the x axis. The algorithms were trained on handwritten 1’s and per method three random 7’s are visualised.

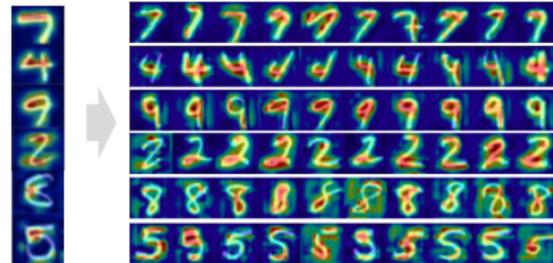


Figure 2: Qualitative results on the MNIST dataset. On the left are the original results and on the right are randomly sampled results produced with method (1) and ReLU.

183 4 Results

184 This section states the results of the experiments described in the previous section.

185 4.1 Attention maps

186 Replicating the results for the MNIST dataset was an involved process – as described in the methodology – leaving us
 187 with various versions of replications. In Figure 1 we present the two different methods based on the pre-trained model
 188 trained on handwritten digit 1’s and tested on handwritten digit 7’s with three different functions (ReLU, absolute value
 189 and sigmoid). Most of these versions would be able to produce the images presented in the paper, except for those
 190 generated using the sigmoid which is why this function was excluded from further experiments. The only difference
 191 between the decent versions is the level of cherry picking Liu et al did. We are assuming a minimal amount, and
 192 therefore conclude that the examples in Figure 4 in the original paper were most likely produced with method (1) and
 193 the ReLU function, even though the paper seems to suggest they were produced using method (2). We chose to proceed
 194 with method (1) as the results resembled those presented in the paper much better. The attention maps for the model
 195 trained on digit 1 and tested on the digits 9, 2 and 4, as well as the model trained on digit 3 and tested on 8 and 5 are
 196 produced with those settings and the pre-trained models. These can be found in Figure 2.

197 We are fairly certain that the differences between the pre-trained model and the models trained from scratch are not due
 198 to differences in architecture or hyperparameters as they were explicitly stated in the addendum paper, which means the
 199 only remaining influential factor is the amount of epochs for which the model was trained. As we reached out to the
 200 authors about the amount of epochs but did not get a response we could not confirm this suspicion.

201 4.2 UCSD Ped1 Results

202 We were able to generate qualitative results which match those presented in the paper to a certain degree, however, they
 203 are not as good in quality. These results can be seen on Figure 3. The frames in the paper could be matched in quality if
 204 we were cherry picking results, which could be one of the explanations for the mismatch, but a more likely explanation
 205 would be that we were simply not able to replicate the results properly.

206 Quantitative results are shown in Table 4. These results show that we were unable to recreate the exact scores presented
 207 by Liu et al. It is worth noting that we were also not able to exactly reproduce the baseline scores presented by Liu et al.
 208 for the VAE model (also *Vanilla-VAE*). The baseline score in the paper reached 0.86 while ours got stuck at 0.82. In
 209 conclusion, all our ROC AUC scores were significantly lower compared to the scores from the paper. The scores were

210 highly dependent on the method, and while method (2) was presented as the anomaly localisation method, method (1)
 211 seems to consistently perform better.

212 An additional topic of uncertainty was how masks without anomalies were handled. As it was not touched upon in the
 213 paper or the addendum paper we explored two options, (1) including all images whether they had anomalies or not
 214 and (2) excluding those frames without any anomalies. Furthermore, because we created this code from scratch we
 215 experimented with various small adjustments to the code, one of which was whether we should upsample or downscale
 216 the images, depending on the convolutional layer. Where the second (downscaling) led to higher scores, we decided to
 217 stick with the first approach as we felt it was more representative of the actual performance of the model (downscaling
 218 increases the probability of getting the right pixel value, so as we apply more convolutions the image gets smaller which
 219 leads to higher scores). Another question left unanswered was whether to do any processing on the generated attention
 220 maps, such as normalization etc. We decided against this, as it did not improve the scores and the normalization in the
 221 code base seemed to only have the goal of displaying the images correctly.

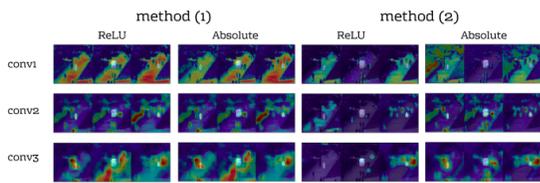


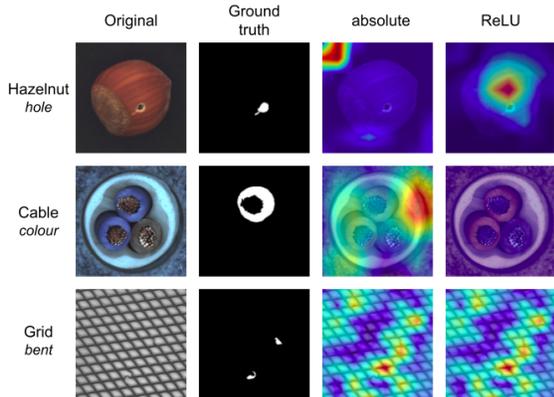
Figure 3: Attention maps for anomaly detection on UCSD, for different layers and generation methods.

	Paper	(1) relu	(1) abs	(2) relu	(2) abs
baseline	0.86	0.82	0.82	0.82	0.82
conv1	0.89	0.71	0.70	0.56	0.59
conv2	0.92	0.80	0.78	0.63	0.69
conv3	0.91	0.68	0.73	0.52	0.64

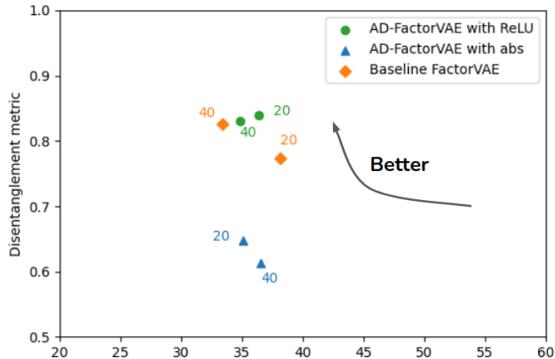
Figure 4: AUROC scores on UCSD

222 4.3 MVTec-AD Results

223 The MVTec-AD results are shown in Table 2. Due to time constraints we were not able to produce scores for all
 224 objects and textures and could only use the (2) second method for attention map generation. We see that our scores
 225 are significantly lower compared to the scores stated in the paper and presented in the second column of Table 2. The
 226 results are also lower compared the scores presented in the original work of Bergmann et al. [12] which are presented
 227 in Table 2 of the original paper by Liu et al. Figure 5a shows some qualitative results from three categories. These
 228 attention maps were one of the best that were produced during our experiments. One can note that these attention maps
 229 are far from accurate in selecting the anomalous regions. Therefore both the quantitative and the qualitative results we
 230 produced are quite different from the results in the original paper. However, in contrast to the UCSD Ped1 experiments,
 231 there is no baseline score we can use to put the scores into perspective.



(a) Qualitative results for Hazelnut, Cable and Grid categories in MVTec-AD. For each category, we show the original image, the ground-truth mask, the attention maps generated with ReLU and with absolute value



(b) Achieved disentanglement scores of FactorVAE and AD-FactorVAE models. The numbers next to the markers indicate γ parameter values.

Figure 5

232 **4.4 Anomaly detection**

233 In conclusion, we can state that we were not able to replicate the results for
 234 anomaly detection proposed by the authors. We therefore cannot support the
 235 claim that the performance of their proposed anomaly localization method
 236 is state-of-the-art by obtaining ROCAUC and IOU scores that are better or
 237 similar to those obtained by the models compared in the original work of
 238 Bergmann et al. Moreover, the attention maps did not exceed the baseline
 239 score for the UCSD Ped1 experiments.

240 There might, however, be room for improvement within our implementations
 241 which could lead to state-of-the-art performance, as stated by Liu et al.
 242 Possibilities for improvement could be training the models for more epochs,
 243 different weight initialization, or changing the way we generate the attention
 244 maps, for which there is some confusion about the usage of method (1) or
 245 method (2), and using ReLU or absolute value.

246 Striking is the fact that method (1) – which was originally just proposed
 247 for visualising the attention maps – consistently outperforms method (2) –
 248 which was specifically proposed for anomaly localization – on the anomaly
 249 localization task and is able to generate results much closer related to those
 250 presented in the paper. The trade-off between the two functions seems to be
 251 quality vs. time. Method (1) has to loop over the entire latent space, making
 252 it much slower, where method (2) just has to propagate the mean vector. We
 253 are unsure which one of these methods was actually used to obtain the final
 254 results presented in the paper.

255 **4.5 Attention Disentanglement Results**

256 As the FactorVAE paper does not mention the weight initialisation used,
 257 we relied on the external code which provides options for both normal and
 258 kaiming normal initialization. However after being unable to reproduce the
 259 baseline results with either of these methods, explicit weight initialization was
 260 removed causing the pytorch default, kaiming uniform, to be applied. This
 261 allowed us to achieve the desired results, but it also shows how sensitive these
 262 methods are to small technical details. Additionally, we had to ignore the
 263 color factor to reach the correct disentanglement scores. Object or background
 264 color cannot be varied randomly in the dSprites dataset hence it always acts
 265 as a fixed factor for the metric.

266 Setups where loss values were aggregated using the mean and a γ value of
 267 20 or 40, converged successfully to the expected reconstruction loss. Results
 268 of their corresponding scores can be seen on Figure 5b. We can clearly see
 269 that using ReLU achieves significantly superior results compared to using
 270 absolute value. However, in general we can see that adding the proposed
 271 loss did not increase the disentanglement score compared to the baseline
 272 FactorVAE.

273 Setups where the loss values were aggregated using a sum or a γ of 100, converged to reconstruction loss values well
 274 over 100. A possible reason to why summing does not work is that then the AD loss gets too large and influences the
 275 overall loss too much. The same could apply to the issues with a large γ parameter and the Total Correlation loss of
 276 FactorVAE, however as this was part of the baseline, this rather suggests that we might have missed key information
 277 about training models with a parameter of this magnitude. Because the authors claim the addition of the AD loss does
 278 not decrease the reconstruction loss, we only determined a model’s disentanglement score if it reached a reconstruction
 279 loss similar to the baseline FactorVAE model.

280 **5 Discussion**

281 In this paper, we have provided insights in the reproducibility attempt of the paper “Towards Visually Explaining
 282 Variational Autoencoders” by Liu et al. Our results show that, with the provided tools and within a time-frame of four
 283 weeks, we have been able to reproduce similar attention maps for the MNIST dataset. The gradient-based attention

Category	theirs	ours ReLU	ours abs
Carpet	0.78 0.1	0.51 0.02	0.47 0.02
Grid	0.73 0.02	0.55 0.02	0.50 0.02
Leather (diff. padding)	0.95 0.24	0.52 0.13	0.67 0.18
Tile	0.80 0.23	0.59 0.14	0.62 0.14
Wood	0.77 0.14	0.42 0.05	0.63 0.07
Bottle	0.87 0.27	0.53 0.08	0.46 0.08
Cable	0.90 0.18	0.50 0.05	0.62 0.08
Capsule	0.74 0.11	0.54 0.03	0.69 0.04
Hazelnut	0.98 0.44	0.63 0.08	0.75 0.09
Metal Nut	0.94 0.49	0.39 0.15	0.46 0.15
Pill	0.83 0.18	0.53 0.05	0.60 0.06
Screw	0.97 0.17	0.52 0.00	0.51 0.00
Toothbrush	0.94 0.14	0.64 0.04	0.65 0.03
Transistor	0.93 0.30	0.47 0.12	0.61 0.15

Table 2: Quantitative results for pixel level segmentation on 14 categories from MVTEC-AD dataset. For each category, we report the area under ROC AUC curve on the top row, and best IOU on the bottom row.

284 maps clearly denoted characteristics that intuitively made sense. We can therefore confirm the first claim by the authors.
285 Besides that, we have not been able to produce any results that were comparable to the ones from paper. This leads to
286 the conclusion that we cannot successfully confirm the second and third claim made by the authors. However, some
287 results and extra experiments indicated that there might be room for improvement within the used methods through
288 hyperparameter tuning or small model changes. Our approach has several weaknesses which will be discussed in the
289 next section.

290 5.1 Strengths and weaknesses

291 Our implementation has several strong points, including:

- 292 • The authors' code was partly available on GitHub containing method (2) for the generation of the attention
293 maps (as explained in Section 3.1) as well as the set-up for the MNIST experiments. This means that for these
294 parts, we are quite certain that our code matches the authors' code.
- 295 • All datasets were publicly available and well documented, causing us to be fairly confident that we used the
296 same data as the authors.
- 297 • For the models that were not implemented in the authors' code, the architectural details were publicly available
298 in the authors' addendum paper. We were therefore able to implement the same models as the authors used.
- 299 • For the FactorVAE model, the authors referred us to a repository which contained an implementation that
300 only needed some small adjustments to give us the correct baseline results. Thus, it seems unlikely that our
301 implementation of the FactorVAE model is incorrect.
- 302 • For the FactorVAE disentanglement metric a Google repository was found that contained a proper implementa-
303 tion, which gives a high probability that our implementation of this metric is correct.
- 304 • In many of our experiments, we have tried several different approaches to obtain the same results as the authors.
305 For example, because we did not obtain correct results for the MNIST experiments, we also implemented the
306 first approach for the attention maps generation even though the paper mentioned the results were obtained
307 using the second approach. For the attention disentanglement we experimented with several attention map
308 generation and loss aggregation methods.

309 However, there are also some weaknesses in our approach:

- 310 • The authors' code and the paper were sometimes inconsistent with each other. Especially for the generation of
311 the attention maps. It was unclear if the results were obtained using ReLU (as stated in the paper) or using
312 absolute value (as stated in the code). We tried to overcome this weakness by implementing both approaches.
313 Furthermore, it was unclear which method to use for which part. The misalignment between the fact that
314 method (2) seems to be presented as the method to use for all experiments in the original paper, while method
315 (1) consistently brings us results that are closer to the results presented is vexing to us and might account for
316 some of the discrepancies between the original results and ours.
- 317 • For each experiment the amount of training epochs was unclear causing us to make guesses due to limited
318 time.
- 319 • It showed to be very difficult to reproduce the qualitative and quantitative results for anomaly localization on
320 both the UCSD Ped1 dataset and the MVTec-AD data even though the models' architectures were specified in
321 the addendum paper. Additional research is needed to discover why our experiments failed to produce the
322 correct scores.
- 323 • It was not possible to implement the exact same data augmentation for the MVTec-AD dataset as the paper did
324 not provide enough information on this topic.
- 325 • The entire implementation for the attention disentanglement was missing in the authors' code. We therefore
326 had to make many guesses about the implementation of this part, which made it difficult to reproduce the
327 qualitative and quantitative results for latent space disentanglement on the dSprites dataset.

328 5.2 Communication with original authors

329 We have communicated with the authors on two occasions. The first time we e-mailed them with a request for the
330 complete code as the public repository had many missing parts. They replied that they were not able to give us
331 their full code as they were planning to patent their implementation in the future. We sent them another e-mail with
332 some follow-up questions about the number of epochs, data-augmentation for the MVTec-AD dataset and method of
333 aggregation of the AD loss. Unfortunately, we did not receive a reply to that last e-mail.

References

- 334
- 335 [1] Wenqian Liu, Runze Li, Meng Zheng, Srikrishna Karanam, Ziyang Wu, Bir Bhanu, Richard J. Radke, and Octavia
336 Camps. Towards visually explaining variational autoencoders, 2020.
- 337 [2] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2014.
- 338 [3] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv
339 Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *International Journal*
340 *of Computer Vision*, 128(2):336–359, Oct 2019.
- 341 [4] Zachary C Lipton. The mythos of model interpretability: In machine learning, the concept of interpretability is
342 both important and slippery. *Queue*, 16(3):31–57, 2018.
- 343 [5] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining the predictions of
344 any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and*
345 *data mining*, pages 1135–1144, 2016.
- 346 [6] Edward Johns, Oisín Mac Aodha, and Gabriel J Brostow. Becoming the expert-interactive multi-class machine
347 teaching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2616–2624,
348 2015.
- 349 [7] Buyu Li, Wanli Ouyang, Lu Sheng, Xingyu Zeng, and Xiaogang Wang. Gs3d: An efficient 3d object detection
350 framework for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
351 *Recognition*, pages 1019–1028, 2019.
- 352 [8] Dakai Jin, Dazhou Guo, Tsung-Ying Ho, Adam P Harrison, Jing Xiao, Chen-Kan Tseng, and Le Lu. Accurate
353 esophageal gross tumor volume segmentation in pet/ct using two-stream chained 3d deep network fusion. In
354 *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 182–191.
355 Springer, 2019.
- 356 [9] Yiming Zuo, Weichao Qiu, Lingxi Xie, Fangwei Zhong, Yizhou Wang, and Alan L Yuille. Craves: Controlling
357 robotic arm with a vision-based economic system. In *Proceedings of the IEEE/CVF Conference on Computer*
358 *Vision and Pattern Recognition*, pages 4214–4223, 2019.
- 359 [10] Li Deng. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE*
360 *Signal Processing Magazine*, 29(6):141–142, 2012.
- 361 [11] Vijay Mahadevan, Wei-Xin LI, Viral Bhalodia, and Nuno Vasconcelos. Anomaly detection in crowded scenes. In
362 *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 1975–1981, 2010.
- 363 [12] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Mvtec ad—a comprehensive real-world
364 dataset for unsupervised anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
365 *and Pattern Recognition*, pages 9592–9600, 2019.
- 366 [13] Hyunjik Kim and Andriy Mnih. Disentangling by factorising, 2019.
- 367 [14] Wonkwang Lee. Pytorch implementation of factorvae proposed in disentangling by factorising, 2018.
- 368 [15] Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dsprites: Disentanglement testing sprites
369 dataset. <https://github.com/deepmind/dsprites-dataset/>, 2017.
- 370 [16] Amir H. Abdi, Purang Abolmaesumi, and Sidney Fels. Variational learning with disentanglement-pytorch. *arXiv*
371 *preprint arXiv:1912.05184*, 2019.