

---

# [Reproducibility Report] Explainable Deep One-Class Classification

---

Anonymous Author(s)

Affiliation

Address

email

## Reproducibility Summary

### 1 **Scope of Reproducibility**

2 Liznerski et al. [23] proposed Fully Convolutional Data Description (FCDD), an explainable version of the Hypersphere  
3 Classifier (HSC) to directly address image anomaly detection (AD) and pixel-wise AD without any post-hoc explainer  
4 methods. The authors claim that FCDD achieves results comparable with the state-of-the-art in sample-wise AD on  
5 Fashion-MNIST and CIFAR-10 and exceeds the state-of-the-art on the pixel-wise task on MVTec-AD. They also give  
6 evidence to show a clear improvement by using few (1 up to 8) real anomalous images in MVTec-AD for supervision at  
7 the pixel level. Finally, a qualitative study with horse images on PASCAL-VOC shows that FCDD can intrinsically  
8 reveal spurious model decisions by providing built-in anomaly score heatmaps.  
9

### 10 **Methodology**

11 We have reproduced the quantitative results in the main text of [23] except for the performance on ImageNet: sample-  
12 wise AD on Fashion-MNIST and CIFAR-10, and pixel-wise AD on MVTec-AD. We used the author's code with GPUs  
13 NVIDIA TITAN X and NVIDIA TITAN Xp. A more detailed look into FCDD's performance variability is presented,  
14 and a Critical Difference (CD) diagram is proposed as a more appropriate tool to compare methods over the datasets in  
15 MVTec-AD. Finally, we study the generalization power of the unsupervised FCDD during training.

### 16 **Results**

17 All per-class performances (in terms of Area Under the ROC Curve (ROC-AUC) [31]) announced in the paper were  
18 replicated with absolute difference of at most 2% and below 1% on average, confirming the paper's claims. We report  
19 the experiments' GPU and CPU memory requirements and their average training time. Our analyses beyond the paper's  
20 scope show that claiming to "exceed the state-of-the-art" should be considered with care, and evidence is given to argue  
21 that the pixel-wise unsupervised FCDD could narrow the gap with its semi-supervised version.

### 22 **What was easy**

23 The paper was clear and explicitly gave many training and hyperparameters details, which were conveniently set as  
24 default in the author's scripts. Their code was well organized and easy to interact with.

### 25 **What was difficult**

26 Using ImageNet proved to be challenging due to its size and need to manually set it up; we could not complete the  
27 experiments on this dataset.

### 28 **Communication with original authors**

29 We reached the main author by e-mail to ask for help with ImageNet and discuss a few practical details. He promptly  
30 replied with useful information.

# 31 1 Introduction

32 Liznerski et al. [23] proposed a deep learning based AD method capable of doing pixel-wise AD (also known as  
33 “anomaly segmentation”) by directly generating anomaly score maps with a loss function based on the Hypersphere  
34 Classifier (HSC) [29], a successor of Deep Support Vector Data Description (DSVDD) [28], using a fully-convolutional  
35 neural network – hence the name Fully Convolutional Data Description (FCDD).

36 By only using convolutions, down-samplings, and batch normalization (no attention mechanism, nor fully connected  
37 layers), an image of dimensions  $C \times H \times W$  (respectively, the number of channels, the height, and the width) is  
38 transformed into a latent representation  $C' \times U \times V$ , where  $U < H$  and  $V < W$ . This low-resolution representation  
39 is a  $U \times V$  grid of  $C'$ -dimensional vectors, from which the pseudo-Huber loss function yields a  $U \times V$  heatmap of  
40 anomaly scores.

41 Each of these  $C'$ -dimensional vectors contains information from a corresponding receptive field within the full resolution  
42 ( $H \times W$ ) image. Evidence [24] suggests that the effective influence of the input pixels decreases Gaussian-ly as their  
43 position is further away from the center of the receptive field. FCDD uses this principle to up-sample the obtained  
44 heatmap back to the original resolution ( $H \times W$ ), therefore directly obtaining a visual, explainable anomaly score map.

45 Finally, FCDD is also adapted to perform anomaly detection at the sample (image) level by taking the average score on  
46 the low resolution anomaly heatmap.

47 **Vocabulary: Sample v.s. Pixel-wise Anomaly Detection** The authors refer to anomaly detection (AD) at the image  
48 level (e.g. given an unseen image, a model trained on horse images should infer if there is a horse present in it or,  
49 otherwise, the image is anomalous) simply as “detection”, while anomaly segmentation/localization (i.e. finding regions,  
50 sets of pixels where there exists anomalous characteristics) is referred as “pixel-wise AD”. Analogously, for the sake of  
51 clarity, we refer to the former as “sample-wise AD”. Both setups are further explained in Section 3.3.

## 52 2 Scope of reproducibility

53 We aimed to reproduce the results announced in [23] to verify the effectiveness of the proposed method both in  
54 sample-wise and pixel-wise anomaly detection. Specifically, we tested the following claims from the original paper:

- 55 1. **Claim 1:** FCDD is comparable with state-of-the-art methods in terms of ROC-AUC in sample-wise anomaly  
56 detection on standard benchmarks (namely, Fashion-MNIST, CIFAR-10, and ImageNet);
- 57 2. **Claim 2:** FCDD exceeds the state-of-the-art on MVTec-AD in anomaly segmentation in the unsupervised  
58 setting in terms of pixel-wise ROC-AUC;
- 59 3. **Claim 3:** FCDD can incorporate real anomalies, and including only a few annotated images ( $\approx 5$ ) containing  
60 real, segmented anomalies, the performance consistently improves;
- 61 4. **Claim 4:** FCDD can reveal spurious model decisions without any extra explanation method on top of it.

62 The experiments on supporting **Claim 1** on Fashion-MNIST and CIFAR-10 have been replicated, as well as all the tests  
63 on MVTec-AD, supporting **Claims 2 and 3**, and the qualitative analysis on PASCAL-VOC, supporting **Claim 4**. We  
64 provide details about computational requirements (CPU memory, GPU memory, and training time) necessary to run  
65 these experiments.

66 **Beyond the paper** Other analyses are proposed on the results obtained from the experiments corresponding to **Claims**  
67 **2 and 3**, which further confirm **Claim 3** but show that **Claim 2** should be taken with consideration. We also investigate  
68 the evolution of the test performance during the optimization in MVTec-AD’s unsupervised setting (see Section 3.3),  
69 revealing opportunity for improvement that could narrow down the gap with the semi-supervised setting.

## 70 3 Methodology

71 We used the author’s code (PyTorch 1.9.1 and Torchvision 0.10.1), publicly available on GitHub [4], to reproduce the  
72 quantitative experiments presented in the main text. It required no external documentation, and the whole reproduction  
73 took roughly one month-person of work.

### 74 3.1 Datasets

75 The proposed method was originally tested [23] on Fashion-MNIST [32], CIFAR-10 [19], ImageNet1k [13], MVTec-  
76 AD [9], and PASCAL VOC [14]. Besides, EMNIST [11], CIFAR-100 [19], and ImageNet21k<sup>1</sup> (version “fall 2011”)  
77 were used as Outlier Exposure (OE) [16] datasets. All the datasets except for ImageNet were publicly available and  
78 automatically downloaded.

79 **ImageNet** We requested access and download ImageNet1k (version “ILSVRC 2012”) from its official website [5].  
80 ImageNet21k (a.k.a. ImageNet22k) was downloaded from `academictorrents.com` [1] because the version used in  
81 the original paper was not available in the official website anymore.

### 82 3.2 Models

83 We used the same neural networks as the original paper, which depend on the dataset:

- 84 • **Fashion-MNIST**: three convolutional layers separated by two max-pool layers, where the first convolution is  
85 followed by a batch normalization and a leaky ReLU;
- 86 • **CIFAR-10**: two convolutions preceded by three blocks, each, composed of a convolution, a batch normaliza-  
87 tion, a leaky ReLU, and a max-pool layer;
- 88 • **MVTec-AD and PASCAL-VOC (Clever Hans)**: the first 10 (frozen) layers from VGG11 pre-trained on  
89 ImageNet followed by two convolutional layers.

### 90 3.3 Experimental setup

91 The paper presents two quantitative experiments: sample-wise (section “4.1 Standard Anomaly Detection Benchmarks”  
92 in [23]) and pixel-wise AD (section “4.2 Explaining Defects in Manufacturing” in [23]), as well as a qualitative  
93 experiment.

94 We followed the same experimental procedure used in [23]: each experiment – i.e. given a dataset, its OE when  
95 applicable, a normal class, and all hyperparameters – was repeated five times, and the reported values are the average  
96 over them unless stated otherwise (e.g. Figure 1).

97 **Sample-wise** Standard one-vs-rest setup, where one class of the given database is chosen as normal and all the others  
98 are used as anomalous. Each image has a binary ground truth signal – logically derived from its label – and the model  
99 assigns an anomaly score to it (therefore “sample-wise”). The metric used is the ROC-AUC on the test split, and all the  
100 classes are evaluated as normal. The datasets used in this experiment and their respective OE dataset is summarized in  
101 Table 1, and its results support **Claim 1**.

Table 1: Sample-wise experiments: tested datasets and their respective OE sources. From the dataset in the column  
"one-vs-rest", one class is used as normal at training and test time while all others are considered anomalies at test time  
only. The column Outlier Exposure (OE) is the dataset used as a source of anomalies at training time. "Experiment  
reference" will further be used to reference these configurations.

One-vs-rest dataset	OE dataset	Experiment reference
Fashion-MNIST	EMNIST	F-MNIST (OE-EMNIST)
	CIFAR-100	F-MNIST (OE-CIFAR-100)
CIFAR-10	CIFAR-100	CIFAR-10
ImageNet1k	ImageNet21k	ImageNet

102 **Pixel-wise** Anomalies are defined at the pixel level (binary segmentation mask where “1” means “anomalous” and “0”  
103 means “normal”) and an image is considered anomalous if it *contains* anomalous pixels although normal pixels are also  
104 present in the image. In each experiment a single class in MVTec-AD is fixed, its normal images are used both for  
105 training and test, and anomalous ones are used for test. As for the anomalous samples at training time, two settings  
106 were tested:

<sup>1</sup>Also known as “ImageNet22k” or “full ImageNet”.

Table 2: Memory requirements and training time using NVIDIA GPUs TITAN X and TITAN Xp (one at a time, indistinctly).

Experiment	CPU memory (Gb)	GPU memory (Gb)	Training duration
F-MNIST (OE-CIFAR-100)	2	1.3	12 min
CIFAR-10	3	1.9	34 min
MVTec-AD unsupervised	38	5.5	1h 13 min
MVTec-AD semi-supervised	33	5.5	41 min
PASCAL VOC (Clever Hans)	5	11.8	21 min

- 107 • **Unsupervised:** synthetic random anomalies are generated using a “confetti noise” (colored blobs added to the  
108 image);
- 109 • **Semi-supervised:** one image per anomaly group (1 up to 8 types depending on the class) is removed from the  
110 test set and used for training.

111 Neither of these settings require an OE dataset because the anomalous samples are either synthetic or real on images of  
112 the nominal class. The performance metric is the ROC-AUC of the anomaly scores at the pixel level. MVTEC-AD is the  
113 only dataset used in this case, and the results of these experiment support **Claims 2 and 3**.

114 **Clever Hans (PASCAL VOC)** About one fifth of the images in the class “horse” in PASCAL VOC [14] contain  
115 a watermark [20], which may cause models to learn spurious features. This is known as the “Clever Hans” effect  
116 as a reference to Hans, a horse claimed to be capable of performing arithmetic operations while it, in fact, read his  
117 master’s reactions [26] – analogously, a model making decisions based on the watermarks would be “cheating” the real  
118 problem. In this experiment, a model is trained using all the classes in PASCAL VOC as normal, only the class “horse”  
119 as anomalous (a swapped one-vs-rest setting), and ImageNet1k as OE dataset. The goal is to qualitatively observe  
120 if one class classifiers are also vulnerable to the Clever Hans effect and show that FCDD transparently reveals such  
121 weaknesses as it intrinsically provides explanations (score heatmaps). This experiment has no quantitative metric but it  
122 supports **Claim 4**.

### 123 3.4 Hyperparameters

124 Running the author’s code with the default parameters, as described in the original paper, did not require any hyperpa-  
125 rameter tuning to achieve the reported results (differences detailed in Section 4) and confirm the authors’ claims. We  
126 underline that the results on MVTEC-AD were obtained using the same hyperparameters on both settings, unsupervised  
127 and semi-supervised.

### 128 3.5 Computational requirements

129 We used the NVIDIA GPUs “TITAN X” [6] and “TITAN Xp” [7] to run our experiments. The two GPUs were used  
130 indistinctly as they have similar characteristics, and only one GPU was used at a time. The GPU and CPU memory  
131 requirements and the average training duration of our experiments are listed in the Table 2 below.

132 CPU memory was recorded with an in-house python script using the library `psutil` [3] at 1Hz. GPU memory was  
133 recorded using `gpustat` [2] at 1Hz. Both memory values are the maximum recorded during the experiments, including  
134 training and inference time. The training duration is an average over all the experiments.

135 On F-MNIST (OE-CIFAR-100) and CIFAR-10, the range of duration did not vary more than two minutes, on MVTEC-  
136 AD unsupervised it ranged from 15 minutes up to one hour, and on MVTEC-AD semi-supervised it ranged from 22  
137 minutes up to one hour and 56 minutes depending on the class.

### 138 3.6 Beyond the paper

139 We propose a more detailed visualization of the distribution of performances (due to random effects, all hyperparameters  
140 held constant) of the two settings (unsupervised and semi-supervised) evaluated on MVTEC-AD, and a critical difference  
141 diagram as alternative evaluation of performance across several datasets (the individual classes in MVTEC-AD).

142 The network architectures used for the experiments on MVTec-AD were pre-trained on ImageNet and most of the  
 143 weights are kept frozen, therefore raising the question of how much of FCDD’s performance is due to the pre-training.  
 144 We took snapshots of the unsupervised model’s weights in order to visualize the evolution of the performance on the  
 145 test set during training.

## 146 4 Results

### 147 4.1 Reproducing the original paper

148 We reproduced the unsupervised and semi-supervised settings for MVTec-AD, and all experiments on Table 1 except  
 149 for ImageNet – due to resource limitations, this experiment could not be completed in time.

150 The results of F-MNIST (OE-EMNIST) were not detailed in the original paper but it is claimed to have a class-mean  
 151 ROC-AUC of  $\sim 3\%$  below F-MNIST (OE-CIFAR-100), and we observed a difference of 2.7%.

152 We summarize the differences between our results and those from the original paper [23] in Table 3. The error margins  
 153 presented are in absolute differences and refer to the ROC-AUC (which is expressed in %) from each individual class’s  
 154 experiment (recall: the mean over five iterations).

Table 3: Differences between the original paper results and ours. All the values are in absolute difference of ROC-AUC, expressed in % (it is **not** a relative error). The columns in "Difference per class" show statistics of the absolute difference of each individual class performance, while the column "Mean ROC-AUC diff." corresponds to the difference measured after the mean is taken over all the classes.

Experiment	ROC-AUC type	N. classes	Diff. per class		Mean ROC-AUC diff.
			Max	Mean	
F-MNIST (OE-CIFAR-100)	sample-wise	10	1%	0.6%	0.01%
CIFAR-10	sample-wise	10	0.5%	0.3%	0.4%
MVTec-AD unsupervised	pixel-wise	15	2%	0.6%	0.2%
MVTec-AD semi-supervised	pixel-wise	15	2%	0.7%	0.4%

155 **Clever Hans (PASCAL VOC)** The experiment on PASCAL VOC (“Clever Hans Effect”) has been manually verified,  
 156 and similar (flawed) explanations on horse images have been observed. Two examples are shown in Figure 5 in  
 157 Appendix A.

### 158 4.2 Beyond the paper

159 Figure 1 further details the performance comparison between the unsupervised and semi-supervised settings on  
 160 MVTec-AD on each class.

161 Figure 2 compares the methods in Table 2 in [23] with a CD diagram using the Wilcoxon-Holm procedure implemented  
 162 by [17]. We replaced the results for FCDD from [23] by our own and copied the others from the literature [10, 30, 25,  
 163 22, 21, 12, 35]. For each class, the methods are sorted by their respective ROC-AUC and assigned a ranking from 1 to  
 164 10 according to their position; then, every pair of methods is compared with the Wilcoxon signed rank test with the  
 165 confidence level  $\alpha = 5\%$ . The CD diagram shows the average ranks of the methods on the horizontal scale, and the red  
 166 bars group methods where each pair of methods are not significantly different according to the Wilcoxon signed rank  
 167 test. The ranks from one to five (six to ten omitted for the sake of brevity) are shown in Table 4 in the Appendix A.

168 Figure 3 shows the test pixel-wise ROC-AUC scores during the optimization of the model used for MVTec-AD with  
 169 the unsupervised setting. Due to time and resources constraints, we ran this experiment on 7 out of the 15 classes in  
 170 MVTec-AD, each of them being evaluated 6 times (a few of them could not finish in time).

Experiments on MVTec-AD: unsupervised V.S. semi-supervised

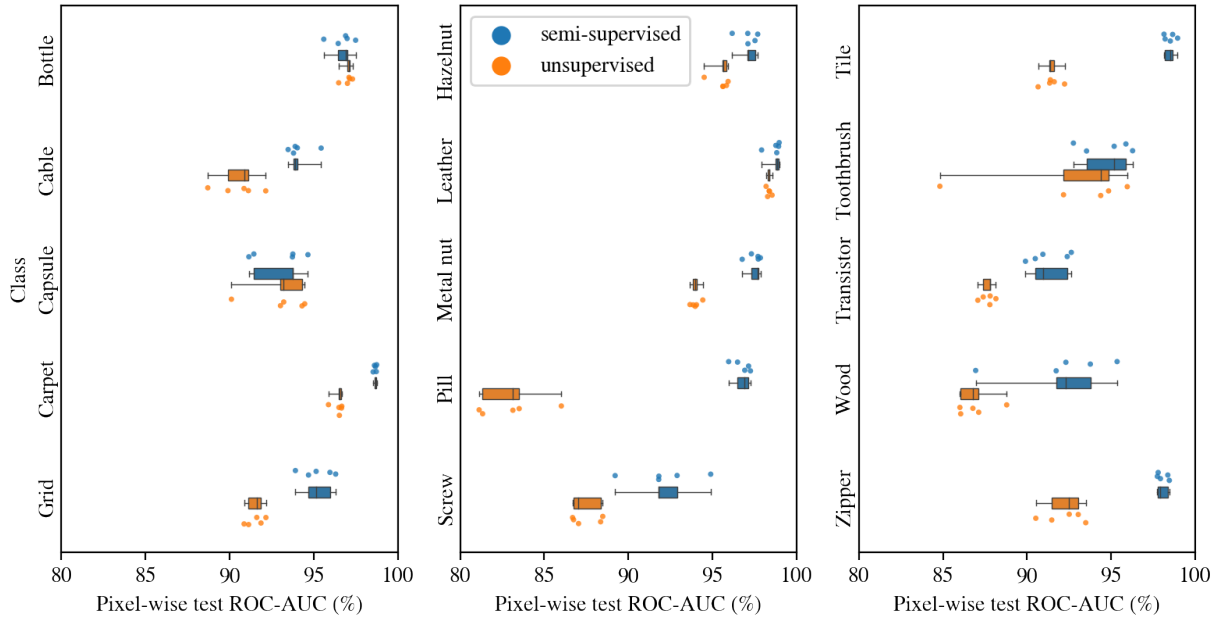


Figure 1: Our experiments on MVTec-AD: unsupervised and semi-supervised settings compared. We display a box plot of the performances (in terms of pixel-wise ROC-AUC on the test set) achieved in different runs along with their individual performances scattered on the  $x$ -axis.

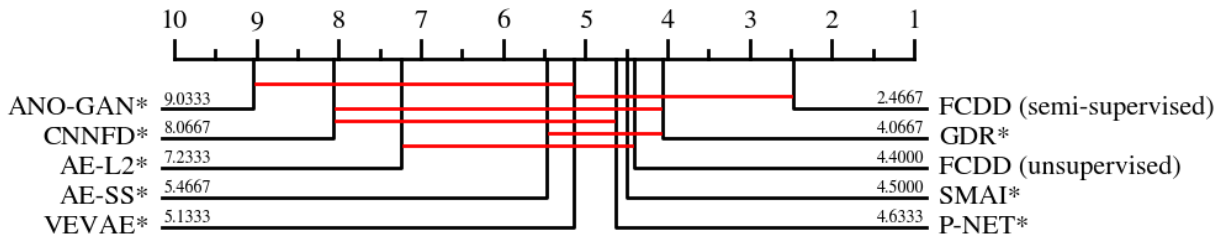


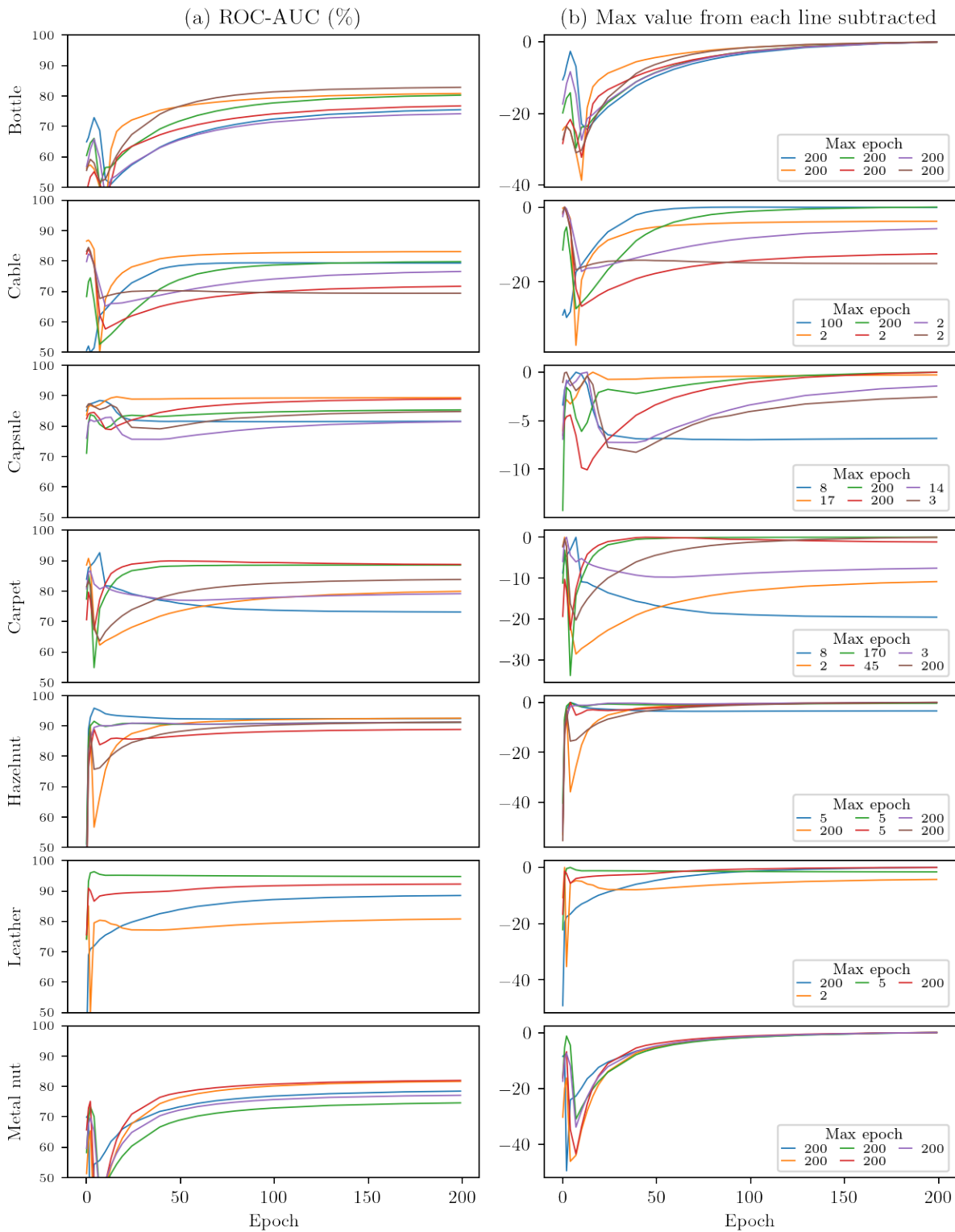
Figure 2: **MVTec-AD Critical Difference diagram.** Using the Table 2 from [23] with the results for FCDD replaced by our own, we build a critical difference diagram using the Wilcoxon-Holm method. Values on the scale are average rankings, and each red line groups a set of methods that are not significantly different in terms of ranking with a confidence level of  $\alpha = 5\%$ . The per-class ROC-AUC values used for FCDD are from our own experiments, and those marked with "\*" were taken from the literature. References: Scores for Self-Similarity (AE-SS) [10], L2 Autoencoder (AE-L2) [10], AnoGAN [30], and CNN Feature Dictionaries (CNNFD) [25] were taken from Table 3 in [10]. Other scores were taken from their respective papers: Visually Explained Variational Autoencoder (VEVAE) from Table 2 in [22], Superpixel Masking and Inpainting (SMAI) from Table 2 in [21], Gradient Descent Reconstruction with VAEs (GDR) from Table 1 in [12], Encoding Structure-Texture Relation with P-Net for AD (P-NET) from Table 6 in [35].

171 **5 Discussion**

172 Our reproduction of the experiments closely agree with the quantitative results published in the original paper. The  
 173 proposed setup is adapted to support the claims announced in the paper hold, and the results corroborate it. We obtained  
 174 results consistently close to the published ones without any further tuning of the parameters or modification of the  
 175 authors' code.

### MVTec-AD (unsupervised) training history

Pixel-level ROC-AUC on the test set measured during the training



Performances were recorded at the following epochs (out of 1 to 200):  
1, 2, 3, 5, 8, 11, 14, 17, 20, 25, 40, 45, 50, 60, 70, 80, 90, 100, 130, 170, 200

Figure 3: MVTec-AD test performance history.

## 176 5.1 What was easy

177 The paper is clear, and it was easy to grasp the core ideas presented in the main text. It also provided enough details  
178 about the experimental setup, including training hyper parameters and network architecture in the appendices.

179 The code was overall well organized and the instructions to use it were direct and easy to use. Conveniently, the  
180 experiments were well encapsulated scripts and default parameters matched those described in the text. In particular,  
181 the experiments are self-documenting (i.e. they keep record of the configurations, logs, and results, etc) and flexible,  
182 allowing the user to change (many) parameters without modifying the code.

## 183 5.2 What was difficult

184 **ImageNet** Using ImageNet was the hardest part. At first, it took about a month to get access to it on the official  
185 website. Then we had to find an alternative source [1] to find the correct version of ImageNet21k (“fall 2011”) because  
186 it was not available on the official website anymore. Basic operations (e.g. decompressing data, moving files) proved  
187 challenging due to its size (1.2 TB compressed), and the instructions to manually prepare this dataset could be more  
188 explicit – we wasted several hours of work because of a few mistakes we made.

189 We could not run the experiments on that dataset with the same hyperparameters because the GPU we dispose of did  
190 not have enough memory (16GB). Although, we note that some solutions like using multiple GPUs or decreasing the  
191 batch size were possible but could not be tried in time.

192 **Minor code issues** There were a few minor bugs, which we corrected without considerable difficulty. They were  
193 mostly related with the script `add_exp_to_base.py`, which automatically configures and launches baseline experi-  
194 ments based on a previously executed one. Finally, the code structure was slightly overcomplicated; e.g. the levels of  
195 abstraction/indirection, specially heritage, could be simpler. Although, we stress that this negative point is minor and  
196 did not cause any critical issues.

## 197 5.3 Communication with original authors

198 We exchanged e-mails with the main author, mostly to ask for help with getting access to the right versions of ImageNet  
199 and executing the experiments on it. He replied promptly, his answers were certainly helpful, and we would like to  
200 express our sincere appreciation.

## 201 5.4 Beyond the paper

202 **MVTec-AD: supervision effect** The visualization proposed in Figure 1 further demonstrates that, with only a few  
203 images of real anomalies added to the training, the model’s performance consistently improves. Only 4/15 classes have  
204 performance overlap, all others show a clear shift in the performance distribution.

205 However, it must be mentioned that the synthetic anomalies ignore the local supervision, therefore making its training  
206 sub-optimal. Figure 4 illustrates this with training images and their respective masks from the class “Pill”: in 4a we see  
207 that the the semi-supervised setting provides pixel-level annotations on the anomalies (the ground truth mask), while in  
208 4b we see that the entire image is considered anomalous in the unsupervised setting. This is a source of sub-optimality  
209 because, in the anomalous images, most pixels are, in fact, normal. In other words, similar images patches, free of  
210 synthetic anomaly, can be found both in normal and anomalous images.

211 Ultimately, this a clear opportunity of improvement that can bring the unsupervised setting’s performance closer to the  
212 semi-supervised setting.

213 **Test performance history** Figure 3 reveals another issue with the method. Take, for instance, the purple and blue  
214 lines in the row “Carpet”; they reach a maximum point at the beginning of the gradient descent then converge to a point  
215 with less and less generalization power. These performance histories are evaluated on the test set, which is assumed to  
216 be unavailable at training time, so this information could not be used to stop the training or reject it. However, this  
217 reveals another opportunity of improvement because the training setting does not push the model to generalize well  
218 enough. Note, in Figure 4b, that the confetti noise also “stains” the background, creating synthetic anomalies “out of  
219 context”, so using more realistic ones could be a solution?



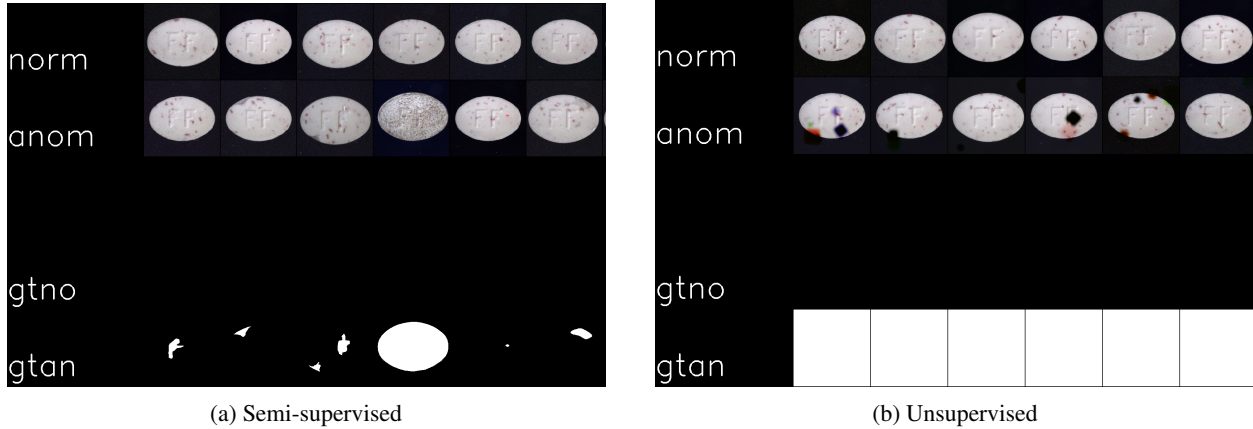


Figure 4: MVTec-AD training images: unsupervised vs. semi-supervised.

220 We also see that some (good) performances are likely due to the pre-training (on ImageNet). For instance, the class  
 221 “Hazelnut” often reaches its maximum test performance (or almost) with few epochs.

222 **Critical Difference (CD) diagram** We propose a CD diagram in Figure 2 as a more appropriate methodology to  
 223 aggregate results of all the classes. As a potential user is choosing an AD method for a new dataset, he or she is looking  
 224 for the method(s) that will most likely be the best on his or her own problem. Therefore, the specific ROC-AUC scores  
 225 on standard datasets have little importance but their relative performances is essential. In other words, what matters for  
 226 a potential user is how the comparison of methods generalizes over specific datasets.

227 The experiments on MVTec-AD do not interact from one class (set as nominal) to another, making them essentially  
 228 independent datasets; therefore, taking the average score over the classes may mislead the analysis. For instance, in  
 229 [23], FCDD unsupervised is claimed to beat the state-of-the-art, although Figure 2 shows that GDR [12] has a better  
 230 average ranking.

231 Note that the red bars in the diagram may give the impression that there is no relevant difference at all; although, it is  
 232 important to observe that it was built considering only 15 datasets (therefore 15 rankings), making the statistical test  
 233 hard, so using more datasets could refine these groups and provide better understanding. Finally, it is worth noting that  
 234 the CD diagram is capable of incorporating new datasets, while the mean score over them would be too much affected  
 235 if some cases are much easier or much harder than the others.

236 **State-of-the-art** It is worth mentioning that more recent methods have claimed better results on the same benchmarks  
 237 used in this work. For instance, at least 5 papers [27, 18, 34, 33, 15] claim to have a mean ROC-AUC above 98% on  
 238 the leaderboard for anomaly segmentation (“pixel-wise AD”) on MVTec-AD in Papers with Code [8]. Unfortunately,  
 239 we did not have the time to fully verify the experimental conditions in the sources but this serves as proxy evidence to  
 240 take these results with consideration.

## 241 References

- 242 [1] Academic torrents, download page for ImageNet22k version “fall 2011”. [https://academictorrents.com/](https://academictorrents.com/details/564a77c1e1119da199ff32622a1609431b9f1c47)  
243 [details/564a77c1e1119da199ff32622a1609431b9f1c47](https://academictorrents.com/details/564a77c1e1119da199ff32622a1609431b9f1c47).
- 244 [2] Documentation page of gpustat. <https://github.com/wookayin/gpustat>.
- 245 [3] Documentation page of the python library psutil. <https://psutil.readthedocs.io/en/latest/>.
- 246 [4] GitHub repository liznerski/fcdd. <https://github.com/liznerski/fcdd>. Downloaded: 2021-12-20. Forked  
247 on commit 7af3d8eadabee81ab8f7db5dea7f8389ef090213.
- 248 [5] ImageNet official page, challenge ILSVRC 2012. [https://image-net.org/challenges/LSVRC/2012/](https://image-net.org/challenges/LSVRC/2012/2012-downloads.php)  
249 [2012-downloads.php](https://image-net.org/challenges/LSVRC/2012/2012-downloads.php).
- 250 [6] NVIDIA’s TITAN X product page. [https://www.nvidia.com/en-us/geforce/products/10series/](https://www.nvidia.com/en-us/geforce/products/10series/titan-x-pascal/)  
251 [titan-x-pascal/](https://www.nvidia.com/en-us/geforce/products/10series/titan-x-pascal/). Accessed: 2022-01-24.
- 252 [7] NVIDIA’s TITAN Xp product page. <https://www.nvidia.com/en-us/titan/titan-xp/>. Accessed: 2022-  
253 01-24.
- 254 [8] Papers with code’s leaderboard on Anomaly Detection on MVTec AD. [https://paperswithcode.com/sota/](https://paperswithcode.com/sota/anomaly-detection-on-mvtec-ad?metric=Segmentation%20AUROC)  
255 [anomaly-detection-on-mvtec-ad?metric=Segmentation%20AUROC](https://paperswithcode.com/sota/anomaly-detection-on-mvtec-ad?metric=Segmentation%20AUROC).
- 256 [9] P. Bergmann, K. Batzner, M. Fauser, D. Sattlegger, and C. Steger. The MVTec Anomaly Detection Dataset: A  
257 Comprehensive Real-World Dataset for Unsupervised Anomaly Detection. *International Journal of Computer*  
258 *Vision*, 129(4):1038–1059, Apr. 2021.
- 259 [10] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger. MVTec AD – A Comprehensive Real-World Dataset for  
260 Unsupervised Anomaly Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
261 *Recognition (CVPR)*, June 2019.
- 262 [11] G. Cohen, S. Afshar, J. Tapson, and A. van Schaik. EMNIST: Extending MNIST to handwritten letters. In *2017*  
263 *International Joint Conference on Neural Networks (IJCNN)*, pages 2921–2926, 2017.
- 264 [12] D. Dehaene, O. Frigo, S. Combrexelle, and P. Eline. Iterative energy-based projection on a normal data manifold  
265 for anomaly localization. page 17, 2020.
- 266 [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical image database.  
267 In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.
- 268 [14] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The Pascal Visual Object Classes  
269 (VOC) Challenge. *International Journal of Computer Vision*, 88(2):303–338, June 2010.
- 270 [15] D. Gudovskiy, S. Ishizaka, and K. Kozuka. CFLOW-AD: Real-Time Unsupervised Anomaly Detection with  
271 Localization via Conditional Normalizing Flows. *arXiv:2107.12571 [cs]*, July 2021. arXiv: 2107.12571 version:  
272 1.
- 273 [16] D. Hendrycks, M. Mazeika, and T. Dietterich. Deep Anomaly Detection with Outlier Exposure. In *International*  
274 *Conference on Learning Representations*, 2019.
- 275 [17] H. Ismail Fawaz, G. Forestier, J. Weber, L. Idoumghar, and P.-A. Muller. Deep learning for time series classification:  
276 a review. *Data Mining and Knowledge Discovery*, 33(4):917–963, 2019.
- 277 [18] J.-H. Kim, D.-H. Kim, S. Yi, and T. Lee. Semi-orthogonal Embedding for Efficient Unsupervised Anomaly  
278 Segmentation. *arXiv:2105.14737 [cs]*, May 2021. arXiv: 2105.14737 version: 1.
- 279 [19] A. Krizhevsky. Learning Multiple Layers of Features from Tiny Images, 2009. Technical report.
- 280 [20] S. Lapuschkin, A. Binder, G. Montavon, K.-R. Muller, and W. Samek. Analyzing Classifiers: Fisher Vectors and  
281 Deep Neural Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*  
282 *(CVPR)*, June 2016.

- 283 [21] Z. Li, N. Li, K. Jiang, Z. Ma, X. Wei, X. Hong, and Y. Gong. Superpixel Masking and Inpainting for Self-  
284 Supervised Anomaly Detection. page 12, 2020.
- 285 [22] W. Liu, R. Li, M. Zheng, S. Karanam, Z. Wu, B. Bhanu, R. J. Radke, and O. Camps. Towards Visually Explaining  
286 Variational Autoencoders. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June  
287 2020.
- 288 [23] P. Liznerski, L. Ruff, R. A. Vandermeulen, B. J. Franks, M. Kloft, and K. R. Müller. Explainable Deep One-Class  
289 Classification. In *International Conference on Learning Representations*, 2021.
- 290 [24] W. Luo, Y. Li, R. Urtasun, and R. Zemel. Understanding the Effective Receptive Field in Deep Convolutional  
291 Neural Networks. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*,  
292 NIPS’16, pages 4905–4913, Red Hook, NY, USA, 2016. Curran Associates Inc. event-place: Barcelona, Spain.
- 293 [25] P. Napoletano, F. Piccoli, and R. Schettini. Anomaly Detection in Nanofibrous Materials by CNN-Based Self-  
294 Similarity. *Sensors (Basel, Switzerland)*, 18(1):209, Jan. 2018.
- 295 [26] O. Pfungst. *Clever Hans (The Horse of Mr. Von Osten) A contribution to experimental animal and human*  
296 *psychology*. Oct. 2010.
- 297 [27] K. Roth, L. Pemula, J. Zepeda, B. Schölkopf, T. Brox, and P. Gehler. Towards Total Recall in Industrial Anomaly  
298 Detection. *arXiv:2106.08265 [cs]*, June 2021. arXiv: 2106.08265 version: 1.
- 299 [28] L. Ruff, R. Vandermeulen, N. Goernitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, and M. Kloft. Deep  
300 one-class classification. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on*  
301 *Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 4393–4402. PMLR, 10–15  
302 Jul 2018.
- 303 [29] L. Ruff, R. A. Vandermeulen, B. J. Franks, K.-R. Müller, and M. Kloft. Rethinking Assumptions in Deep Anomaly  
304 Detection. *arXiv:2006.00339 [cs, stat]*, July 2021. arXiv: 2006.00339.
- 305 [30] T. Schlegl, P. Seeböck, S. M. Waldstein, U. Schmidt-Erfurth, and G. Langs. Unsupervised Anomaly Detection  
306 with Generative Adversarial Networks to Guide Marker Discovery. In M. Niethammer, M. Styner, S. Aylward,  
307 H. Zhu, I. Oguz, P.-T. Yap, and D. Shen, editors, *Information Processing in Medical Imaging*, pages 146–157,  
308 Cham, 2017. Springer International Publishing.
- 309 [31] K. A. Spackman. Signal Detection Theory: Valuable Tools for Evaluating Inductive Learning. In *Proceedings of*  
310 *the Sixth International Workshop on Machine Learning*, pages 160–163, Ithaca, New York, USA, 1989. Morgan  
311 Kaufmann Publishers Inc., 340 Pine Street, Sixth Floor, San Francisco, CA, USA.
- 312 [32] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning  
313 Algorithms, 2017. arXiv: 1708.07747.
- 314 [33] J. Yu, Y. Zheng, X. Wang, W. Li, Y. Wu, R. Zhao, and L. Wu. FastFlow: Unsupervised Anomaly Detection and  
315 Localization via 2D Normalizing Flows. *arXiv:2111.07677 [cs]*, Nov. 2021. arXiv: 2111.07677 version: 2.
- 316 [34] Y. Zheng, X. Wang, R. Deng, T. Bao, R. Zhao, and L. Wu. Focus Your Distribution: Coarse-to-Fine Non-  
317 Contrastive Learning for Anomaly Detection and Localization. *arXiv:2110.04538 [cs]*, Oct. 2021. arXiv:  
318 2110.04538 version: 1.
- 319 [35] K. Zhou, Y. Xiao, J. Yang, J. Cheng, W. Liu, W. Luo, Z. Gu, J. Liu, and S. Gao. Encoding Structure-Texture  
320 Relation with P-Net for Anomaly Detection in Retinal Images. In A. Vedaldi, H. Bischof, T. Brox, and J.-M.  
321 Frahm, editors, *Computer Vision – ECCV 2020*, pages 360–377, Cham, 2020. Springer International Publishing.

## 322 A Supplementary details

Normal Class	5	4	3	2	1
Bottle	GDR* 92.0	AE-SS* 93.0	<b>FCDD (SS)</b> 96.7	<b>FCDD (U)</b> 97.0	P-NET* 99.0
Cable	VEVAE* 90.0	<b>FCDD (U)</b> 90.5	GDR* 91.0	SMAI* 92.0	<b>FCDD (SS)</b> 94.1
Capsule	GDR* 92.0	<b>FCDD (SS)</b> 92.9	SMAI* 93.0	<b>FCDD (U)</b> 93.0	AE-SS* 94.0
Carpet	VEVAE* 78.0	AE-SS* 87.0	SMAI* 88.0	<b>FCDD (U)</b> 96.4	<b>FCDD (SS)</b> 98.7
Grid	AE-SS* 94.0	<b>FCDD (SS)</b> 95.2	GDR* 96.0	SMAI* 97.0	P-NET* 98.0
Hazelnut	P-NET* 97.0	SMAI* 97.0	<b>FCDD (SS)</b> 97.1	GDR* 98.0	VEVAE* 98.0
Leather	P-NET* 89.0	GDR* 93.0	VEVAE* 95.0	<b>FCDD (U)</b> 98.4	<b>FCDD (SS)</b> 98.7
Metal nut	GDR* 91.0	SMAI* 92.0	<b>FCDD (U)</b> 94.0	VEVAE* 94.0	<b>FCDD (SS)</b> 97.5
Pill	AE-SS* 91.0	P-NET* 91.0	SMAI* 92.0	GDR* 93.0	<b>FCDD (SS)</b> 96.8
Screw	AE-L2* 96.0	AE-SS* 96.0	SMAI* 96.0	VEVAE* 97.0	P-NET* 100.0
Tile	VEVAE* 80.0	<b>FCDD (U)</b> 91.4	CNNFD* 93.0	P-NET* 97.0	<b>FCDD (SS)</b> 98.5
Toothbrush	VEVAE* 94.0	<b>FCDD (SS)</b> 94.7	SMAI* 96.0	GDR* 99.0	P-NET* 99.0
Transistor	<b>FCDD (U)</b> 87.6	AE-SS* 90.0	<b>FCDD (SS)</b> 91.3	GDR* 92.0	VEVAE* 93.0
Wood	GDR* 84.0	<b>FCDD (U)</b> 86.9	CNNFD* 91.0	<b>FCDD (SS)</b> 92.0	P-NET* 98.0
Zipper	AE-SS* 88.0	P-NET* 90.0	SMAI* 90.0	<b>FCDD (U)</b> 92.2	<b>FCDD (SS)</b> 98.1

Table 4: Method rankings on MVTEC-AD based on pixel-wise ROC-AUC. Using the Table 2 from [23], we compare the methods by normal class individually sorting the performances by pixel-wise ROC-AUC. The numbers in column names indicate the ranking (from 1 to 10); only the first 5 are displayed for the sake brevity. FCDD "unsupervised" and "semi-supervised" versions are respectively indicated by "(U)" and "(SS)", and their original values have been replaced by our own experiments' results.

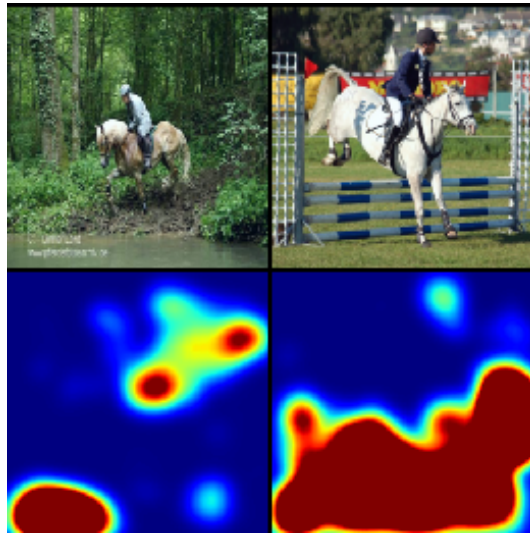


Figure 5: Heatmaps from the experiment on PASCAL-VOC where the Clever Hans effect can be observed.