
A Unified Model for Multi-class Anomaly Detection

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Despite the rapid advance of unsupervised anomaly detection, existing methods
2 require to train separate models for different objects. In this work, we present
3 *UniAD* that accomplishes anomaly detection for multiple classes with a unified
4 framework. Under such a challenging setting, popular reconstruction networks
5 may fall into an “identical shortcut”, where both normal and anomalous samples
6 can be well recovered, and hence fail to spot outliers. To tackle this obstacle, we
7 make three improvements. First, we revisit the formulations of fully-connected
8 layer, convolutional layer, as well as attention layer, and confirm the important role
9 of query embedding (*i.e.*, within attention layer) in preventing the network from
10 learning the shortcut. We therefore come up with a *layer-wise query decoder* to
11 help model the multi-class distribution. Second, we employ a *neighbor masked*
12 *attention* module to further avoid the information leak from the input feature to
13 the reconstructed output feature. Third, we propose a *feature jittering* strategy
14 that urges the model to recover the correct message even with noisy inputs. We
15 evaluate our algorithm on MVTec-AD and CIFAR-10 datasets, where we surpass
16 the state-of-the-art alternatives by a sufficiently large margin. For example, when
17 learning a unified model for 15 categories in MVTec-AD, we surpass the second
18 competitor on the tasks of both anomaly detection (from 88.1% to 96.5%) and
19 anomaly localization (from 89.5% to 96.8%). Code will be made publicly available.

20 1 Introduction

21 Anomaly detection has found an increasingly wide utilization in manufacturing defect detection [4],
22 medical image analysis [16], and video surveillance [45]. Considering the highly diverse anomaly
23 types, a common solution is to model the distribution of normal samples and then identify anomalous
24 ones via finding outliers. It is therefore crucial to learn a compact boundary for normal data, as shown
25 in Fig. 1a. For this purpose, existing methods [6, 10, 24, 26, 47, 50] propose to train separate models
26 for different classes of objects, like in Fig. 1c. However, such a one-class-one-model scheme could
27 be memory-consuming especially along with the number of classes increasing, and also uncongenial
28 to the scenarios where the normal samples manifest themselves in a large intra-class diversity (*i.e.*,
29 one object consists of various types).

30 In this work, we target a more practical task, which is to detect anomalies from different object
31 classes with a unified framework. The task setting is illustrated in Fig. 1d, where the training data
32 covers normal samples from a range of categories, and the learned model is asked to accomplish
33 anomaly detection for all these categories without any fine-tuning. It is noteworthy that the categorical
34 information (*i.e.*, class label) is inaccessible at both the training and the inference stages, considerably
35 easing the difficulty of data preparation. Nonetheless, solving such a task is fairly challenging. Recall
36 that the rationale behind unsupervised anomaly detection is to model the distribution of normal data
37 and find a compact decision boundary as in Fig. 1a. When it comes to the multi-class case, we expect
38 the model to capture the distribution of all classes simultaneously such that they can share the same

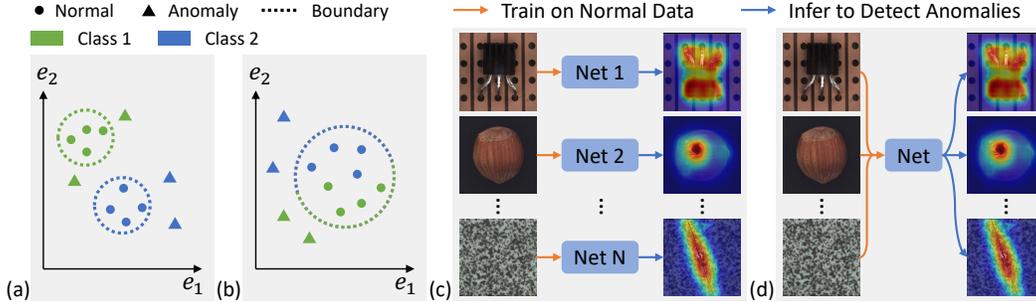


Figure 1: **Task setting of unified anomaly detection.** (a) Existing methods learn separate decision boundaries for different object classes, while (b) our approach models the multi-class data distribution such that one boundary is enough to spot outliers regarding all categories. As a result, we escape from the conventional one-class-one-model paradigm in (c), and manage to accomplish anomaly detection for various classes with a unified framework in (d).

39 boundary as in Fig. 1b. But if we focus on a particular category, say the green one in Fig. 1b, all
 40 the samples from other categories should be considered as anomalies no matter whether they are
 41 normal (*i.e.*, blue circles) or anomalous (*i.e.*, blue triangles) themselves. From this perspective, how
 42 to accurately model the multi-class distribution becomes vital.

43 A widely used approach to learning the normal data distribution draws support from image (or feature)
 44 reconstruction [2, 5, 25, 38, 49], which assumes that a well-trained model always produces normal
 45 samples regardless of the defects within the inputs. In this way, there will be large reconstruction
 46 errors for anomalous samples, making them distinguishable from the normal ones. However, we
 47 find that popular reconstruction networks suggest unsatisfying performance on the challenging task
 48 studied in this work. They typically fall into an “identity shortcut”, which appears as returning a
 49 direct copy of the input disregarding its content.¹ As a result, even anomalous samples can be well
 50 recovered with the learned model and hence become hard to detect.

51 To address this issue, we carefully tailor a feature reconstruction framework that prevents the model
 52 from learning the shortcut. First, we revisit the formulations of fully-connected layer, convolutional
 53 layer, as well as attention layer used in neural networks, and observe that both fully-connected
 54 layer and convolutional layer face the risk of learning a trivial solution. This drawback is further
 55 amplified under the multi-class setting in that the normal data distribution becomes far more complex.
 56 Instead, the attention layer is sheltered from such a risk, benefiting from a learnable query embedding
 57 (see Sec. 3.1). Accordingly, we propose a *layer-wise query decoder* to intensify the use of query
 58 embedding. Second, we argue that the full attention (*i.e.*, every feature point relates to each other)
 59 also contributes to the shortcut issue, because it offers the chance of directly copying the input to
 60 the output. To avoid the information leak, we employ a *neighbor masked attention* module, where a
 61 feature point relates to neither itself nor its neighbors. Third, inspired by Bengio et al. [3], we propose
 62 a *feature jittering* strategy, which requires the model to recover the source message even with noisy
 63 inputs. All these designs help the model escape from the “identity shortcut”, as shown in Fig. 2b.
 64 Extensive experiments on MVTec-AD [4] and CIFAR-10 [22] demonstrate the sufficient superiority
 65 of our approach, which we call *UniAD*, over existing alternatives under the unified task setting. For
 66 instance, when learning a single model for 15 categories in MVTec-AD, we achieve state-of-the-art
 67 performance on the tasks of both anomaly detection and anomaly localization, boosting the AUROC
 68 from 88.1% to 96.5% and from 89.5% to 96.8%, respectively.

69 2 Related work

70 **Anomaly detection.** 1) *Classical approaches* extend classical machine learning methods for one-class
 71 classification, such as one-class support vector machine (OC-SVM) [37] and support vector data
 72 description (SVDD) [34, 40]. Patch-level embedding [47], geometric transformation [17], and elastic
 73 weight consolidation [32] are incorporated for improvement. 2) *Pseudo-anomaly* converts anomaly
 74 detection to supervised learning, including classification [24, 31, 44], image denoising [50], and hyper-

¹A detailed analysis can be found in Sec. 3.1 and Fig. 2.

75 sphere segmentation [26]. However, these methods partly rely on how well proxy anomalies match
 76 real anomalies that are not known [12]. 3) *Modeling then comparison* assumes that the pre-trained
 77 network is capable of extracting discriminative features for anomaly detection [10, 33]. PaDiM [10]
 78 and MDND [33] extract pre-trained features to model normal distribution, then utilize a distance
 79 metric to measure the anomalies. Nevertheless, these methods need to memorize and model all normal
 80 features, thus are computationally expensive. 4) *Knowledge distillation* proposes that the student
 81 distilled by a teacher on normal samples could only extract normal features [6, 12, 36, 43, 44]. Recent
 82 works mainly focus on model ensemble [6], feature pyramid [36, 43], and reverse distillation [12].

83 **Reconstruction-based anomaly detection.** These methods rely on the hypothesis that reconstruction
 84 models trained on normal samples only succeed in normal regions, but fail in anomalous regions [5,
 85 25, 35]. Early attempts include Auto-Encoder (AE) [5, 8], Variational Auto-Encoder (VAE) [21, 25],
 86 and Generative Adversarial Net (GAN) [2, 29, 35, 49]. However, these methods face the problem that
 87 the model could learn tricks that the anomalies are also restored well. Accordingly, researchers adopt
 88 different strategies to tackle this issue, such as adding instructional information (*i.e.*, structural [51]
 89 or semantic [38, 45]), memory mechanism [18, 19, 28], iteration mechanism [11], image masking
 90 strategy [46], and pseudo-anomaly [8, 31]. Recently, DRAEM [50] first recovers the pseudo-anomaly
 91 disturbed normal images for representation, then utilizes a discriminative net to distinguish the
 92 anomalies, achieving excellent performance. However, DRAEM [50] ceases to be effective under
 93 the unified case. Moreover, there is still an important aspect that has not been well studied, *i.e.*,
 94 what architecture is the best reconstruction model? In this paper, we first compare and analyze three
 95 popular architectures including MLP, CNN, and transformer. Then, accordingly, we base on the
 96 transformer and further design three improvements, which compose our UniAD.

97 **Transformer in anomaly detection.** Transformer [41] with attention mechanism, first proposed in
 98 natural language processing, has been successfully used in computer vision [7, 15]. Some attempts
 99 try to utilize transformer for anomaly detection. InTra [30] adopts transformer to recover the image
 100 by recovering all masked patches one by one. VT-ADL [27] and AnoVit [48] both apply transformer
 101 encoder to reconstruct images. However, these methods directly utilize vanilla transformer, and do
 102 not figure out why transformer brings improvement. In contrast, we confirm the efficacy of the query
 103 embedding to prevent the shortcut, and accordingly design a layer-wise query decoder. Also, to avoid
 104 the information leak of the full attention, we employ a neighbor masked attention module.

105 3 Method

106 3.1 Revisiting feature reconstruction for anomaly detection

107 In Fig. 2, following the feature reconstruction paradigm [38], we build an MLP, a CNN, and a
 108 transformer (with query embedding) to reconstruct the features extracted by a pre-trained backbone.
 109 The reconstruction errors represent the anomaly possibility. The architectures of the three networks
 110 are given in *Supplementary Material*. The metric is evaluated every 10 epochs. Note that the periodic
 111 evaluation is *impractical* since anomalies are not available during training. As shown in Fig. 2a, after
 112 a period of training, the performances of the three networks decrease severely with the losses going
 113 extremely small. We attribute this to the problem of “identical shortcut”, where both normal and
 114 anomalous regions can be well recovered, thus failing to spot anomalies. This speculation is verified
 115 by the visualization results in Fig. 2b (more results in *Supplementary Material*). However, compared
 116 with MLP and CNN, the transformer suffers from a much smaller performance drop, indicating a
 117 slighter shortcut problem. This encourages us to analyze as follows.

118 We denote the features in a normal image as $\mathbf{x}^+ \in \mathbb{R}^{K \times C}$, where K is the feature number, C is
 119 the channel dimension. The batch dimension is omitted for simplicity. Similarly, the features in an
 120 anomalous image are denoted as $\mathbf{x}^- \in \mathbb{R}^{K \times C}$. The reconstruction loss is chosen as the MSE loss.
 121 We provide a rough analysis using a simple 1-layer network as the reconstruction net, which is trained
 122 with \mathbf{x}^+ and tested to detect anomalous regions in \mathbf{x}^- .

123 **Fully-connected layer in MLP.** Denote the weights and bias in this layer as $\mathbf{w} \in \mathbb{R}^{C \times C}$, $\mathbf{b} \in \mathbb{R}^C$,
 124 respectively, this layer can be represented as,

$$\mathbf{y} = \mathbf{x}^+ \mathbf{w} + \mathbf{b} \in \mathbb{R}^{K \times C}. \quad (1)$$

125 With the MSE loss pushing \mathbf{y} to \mathbf{x}^+ , the model may take shortcut to regress $\mathbf{w} \rightarrow \mathbf{I}$ (identity matrix),
 126 $\mathbf{b} \rightarrow \mathbf{0}$. Ultimately, this model could also reconstruct \mathbf{x}^- well, failing in anomaly detection.

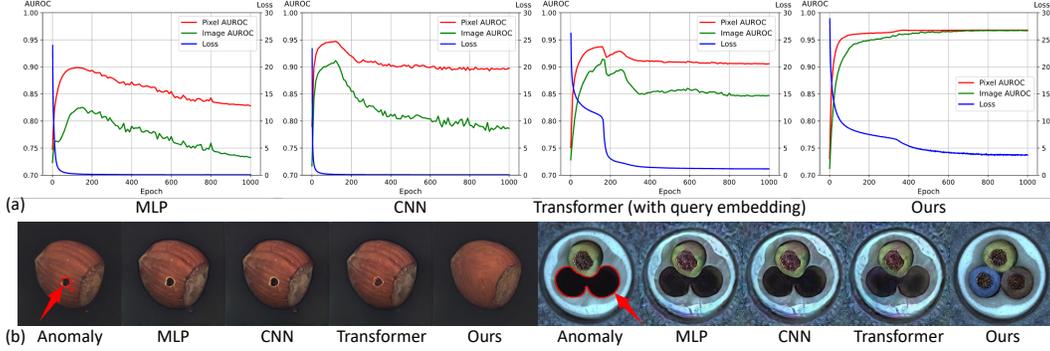


Figure 2: **Comparison among MLP, CNN, transformer, and our UniAD on MVtec-AD [4].** (a) Training loss (blue) as well as the testing AUROC on anomaly detection (green) and localization (red). During the training of MLP, CNN, and transformer, the reconstruction error keeps going smaller on normal samples, but the performance on anomalies suffers from a severe drop after reaching the peak. This is caused by the model learning an “identical shortcut”, which tends to directly copy the input as the output regardless of whether it is normal or anomalous. (b) Visual explanation of the shortcut issue, where the anomalous samples can be well recovered and hence become hard to detect from normal ones. In contrast, UniAD overcomes such a problem and manages to *reconstruct anomalies as normal samples*. It is noteworthy that all models are learned for feature reconstruction and a separate decoder is employed to render images from features. This decoder is *only* used for visualization.

127 **Convolutional layer in CNN.** A convolutional layer with 1×1 kernel is equivalent to a fully-
 128 connected layer. Besides, An $n \times n$ ($n > 1$) kernel has more parameters and larger capacity, and can
 129 complete whatever 1×1 kernel can. Thus, this layer also has the chance to learn a shortcut.

130 **Transformer with query embedding.** In such a model, there is an attention layer with a learnable
 131 query embedding, $\mathbf{q} \in \mathbb{R}^{K \times C}$. When using this layer as the reconstruction model, it is denoted as,

$$\mathbf{y} = \text{softmax}(\mathbf{q}(\mathbf{x}^+)^T / \sqrt{C})\mathbf{x}^+ \in \mathbb{R}^{K \times C}. \quad (2)$$

132 To push \mathbf{y} to \mathbf{x}^+ , the attention map, $\text{softmax}(\mathbf{q}(\mathbf{x}^+)^T / \sqrt{C})$, should approximate \mathbf{I} (identity matrix),
 133 so \mathbf{q} must be highly related to \mathbf{x}^+ . Considering that \mathbf{q} in the trained model is relevant to normal
 134 samples, the model could not reconstruct \mathbf{x}^- well. The ablation study in Sec. 4.5 shows that without
 135 the query embedding, the performance of transformer drops dramatically by 13.4% and 18.1% in
 136 pixel ROAUC and image ROAUC, respectively. Thus the query embedding is of vital significance to
 137 model the normal distribution.

138 However, transformer still suffers from the shortcut problem, which inspires our three improvements.
 139 1) According to that the query embedding can prevent reconstructing anomalies, we design a Layer-
 140 wise Query Decoder (LQD) by adding the query embedding in each decoder layer rather than only
 141 the first layer in vanilla transformer. 2) We suspect that the full attention increases the possibility
 142 of the shortcut. Since one token could see itself and its neighbor regions, it is easy to reconstruct
 143 by simply copying. Thus we mask the neighbor tokens when calculating the attention map, called
 144 Neighbor Masked Attention (NMA). 3) We employ a Feature Jittering (FJ) strategy to disturb the
 145 input features, leading the model to learn normal distribution from denoising. Benefiting from these
 146 designs, our UniAD achieves satisfying performance, as illustrated in Fig. 2.

147 3.2 Improving feature reconstruction for unified anomaly detection

148 **Overview.** As shown in Fig. 3, our UniAD is composed of a Neighbor Masked Encoder (NME)
 149 and a Layer-wise Query Decoder (LQD). Firstly, the feature tokens extracted by a fixed pre-trained
 150 backbone are further integrated by NME to derive the encoder embeddings. Then, in each layer
 151 of LQD, a learnable query embedding is successively fused with the encoder embeddings and the
 152 outputs of the previous layer (self-fusion for the first layer). The feature fusion is completed by
 153 the Neighbor Masked Attention (NMA). The final outputs of LQD are viewed as the reconstructed
 154 features. Also, we propose a Feature Jittering (FJ) strategy to add perturbations to the input features,
 155 leading the model to learn normal distribution from the denoising task. Finally, the results of anomaly
 156 localization and detection are obtained through the reconstruction differences.

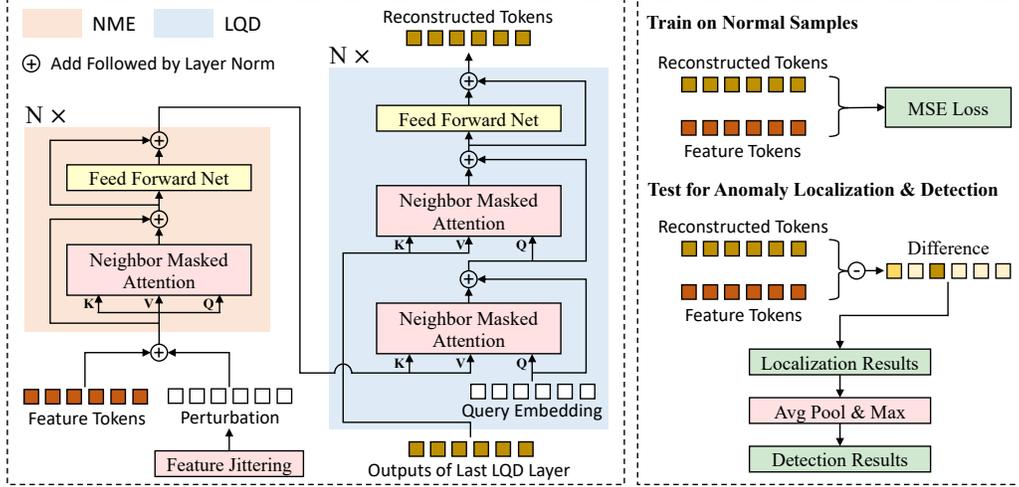


Figure 3: **Framework** of UniAD, consisting of a Neighbor Masked Encoder (NME) and a Layer-wise Query Decoder (LQD). Each layer in LQD employs a *learnable query embedding* to help model the complex training data distribution. The full attention in transformer is replaced by *neighbor masked attention* to avoid the information leak from the input to the output. The *feature jittering* strategy encourages the model to recover the correct message with noisy inputs. All the three improvements assist the model against learning the “identical shortcut” (see Sec. 3.1 and Fig. 2 for details).

157 **Neighbor masked attention.** We suspect
 158 that the full attention in vanilla trans-
 159 former [41] contributes to the “identical
 160 shortcut”. In full attention, one token is
 161 permitted to see itself, so it will be easy to
 162 reconstruct by simply copying. Moreover,
 163 considering that the feature tokens are
 164 extracted by a CNN backbone, the neigh-
 165 bor tokens must share lots of similarities.
 166 Therefore, we propose to mask the neigh-
 167 bor tokens when calculating the attention map,
 168 called Neighbor Masked Attention (NMA).
 169 Note that the neighbor region is defined in
 170 the 2D space, as shown in Fig. 4.

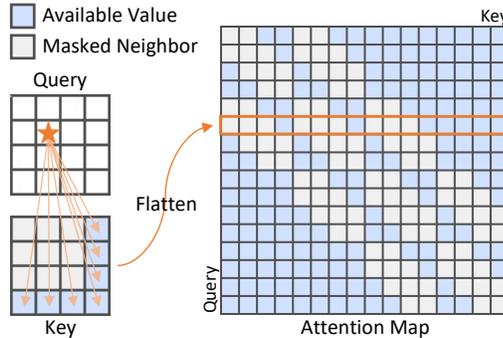


Figure 4: **Illustration of neighbor masked attention**, where pixels relate to neither themselves nor neighbors.

171 **Neighbor masked encoder.** The encoder follows the standard architecture in vanilla transformer.
 172 Each layer consists of an attention module and a Feed-Forward Network (FFN). However, the full
 173 attention is replaced by our proposed NMA to prevent the information leak.

174 **Layer-wise query decoder.** It is analyzed in Sec. 3.1 that the query embedding could help prevent
 175 reconstructing anomalies well. However, there is only one query embedding in the vanilla transformer.
 176 Therefore, we design a Layer-wise Query Decoder (LQD) to intensify the use of query embedding,
 177 as shown in Fig. 3. Specifically, in each layer of LQD, a learnable query embedding is first fused
 178 with the encoder embeddings, then integrated with the outputs of the previous layer (self-integration
 179 for the first layer). The feature fusion is implemented by NMA. Following the vanilla transformer, a
 180 2-layer FFN is applied to handle these fused tokens, and the residual connection is utilized to facilitate
 181 the training. The final outputs of LQD serve as the reconstructed features.

182 **Feature jittering.** Inspired by Denoising Auto-Encoder (DAE) [3, 42], we add perturbations to feature
 183 tokens, guiding the model to learn knowledge of normal samples by the denoising task. Specifically,
 184 for a feature token, $f_{tok} \in \mathbb{R}^C$, we sample the disturbance D from a Gaussian distribution,

$$D \sim N(\mu = 0, \sigma^2 = (\alpha \frac{\|f_{tok}\|_2}{C})^2), \quad (3)$$

185 where α is the jittering scale to control the noisy degree. Also, the sampled disturbance is added to
 186 f_{tok} with a fixed jittering probability, p .

187 3.3 Implementation details

188 **Feature extraction.** We adopt a fixed EfficientNet-b4 [39] pre-trained on ImageNet [13] as the
189 feature extractor. The features from stage-1 to stage-4 are selected. Here the stage means the
190 combination of blocks that have the same size of feature maps. Then these features are resized to the
191 same size, and concatenated along channel dimension to form a feature map, $\mathbf{f}_{org} \in \mathbb{R}^{C_{org} \times H \times W}$.

192 **Feature reconstruction.** The feature map, \mathbf{f}_{org} , is first tokenized to $H \times W$ feature tokens, followed
193 by a linear projection to reduce C_{org} to a smaller channel, C . Then these tokens are processed by
194 NME and LQD. The learnable position embeddings [14, 15] are added in attention modules to inform
195 the spatial information. Afterward, another linear projection is used to recover the channel from C to
196 C_{org} . After reshape, the reconstructed feature map, $\mathbf{f}_{rec} \in \mathbb{R}^{C_{org} \times H \times W}$, is finally obtained.

197 **Objective function.** Our model is trained with the MSE loss as,

$$\mathcal{L} = \frac{1}{H \times W} \|\mathbf{f}_{org} - \mathbf{f}_{rec}\|_2^2. \quad (4)$$

198 **Inference for anomaly localization.** The result of anomaly localization is an anomaly score map,
199 which assigns an anomaly score for each pixel. Specifically, the anomaly score map, s , is calculated
200 as the L2 norm of the reconstruction differences as,

$$\mathbf{s} = \|\mathbf{f}_{org} - \mathbf{f}_{rec}\|_2 \in \mathbb{R}^{H \times W}. \quad (5)$$

201 Then s is up-sampled to the image size with bi-linear interpolation to obtain the localization results.

202 **Inference for anomaly detection.** Anomaly detection aims to detect whether an image contains
203 anomalous regions. We transform the anomaly score map, s , to the anomaly score of the image by
204 taking the maximum value of the averagely pooled s .

205 4 Experiment

206 4.1 Datasets and metrics

207 **MVTec-AD** [4] is a comprehensive, multi-object, multi-defect industrial anomaly detection dataset
208 with 15 classes. For each anomalous sample in the test set, the ground-truth includes both image
209 label and anomaly segmentation. In the existing literature, only the separate case is researched. In
210 this paper, we introduce the unified case, where only one model is used to handle all categories.

211 **CIFAR-10** [22] is a classical image classification dataset with 10 categories. Existing methods [6, 23,
212 36] evaluate CIFAR-10 mainly in the *one-versus-many* setting, where one class is viewed as normal
213 samples, and others serve as anomalies. [Semantic AD](#) [1, 9] proposes a *many-versus-one* setting,
214 [treating one class as anomalous and the remaining classes as normal](#). Different from both, we propose
215 a unified case (*many-versus-many* setting), which is detailed in Sec. 4.4.

216 **Metrics.** Following prior works [4, 6, 50], the Area Under the Receiver Operating Curve (AUROC)
217 is used as the evaluation metric for anomaly detection.

218 4.2 Anomaly detection on MVTEC-AD

219 **Setup.** Anomaly detection aims to detect whether an image contains anomalous regions. The
220 performance is evaluated on MVTEC-AD [4]. The image size is selected as 224×224 , and the size for
221 resizing feature maps is set as 14×14 . The feature maps from stage-1 to stage-4 of EfficientNet-b4
222 [39] are resized and concatenated together to form a 272-channel feature map. The reduced channel
223 dimension is set as 256. AdamW optimizer [20] with weight decay 1×10^{-4} is used. Our model is
224 trained for 1000 epochs on 8 GPUs (NVIDIA Tesla V100 16GB) with batch size 64. The learning
225 rate is 1×10^{-4} initially, and dropped by 0.1 after 800 epochs. The neighbor size, jittering scale, and
226 jittering probability are set as 7×7 , 20, and 1, respectively. The evaluation is run with 5 random seeds.
227 [In both the separate case and the unified case, the reconstruction models are trained from the scratch.](#)

228 **Baselines.** Our approach is compared with baselines including: US [6], PSVDD [47], PaDiM [10],
229 CutPaste [24], MKD [36], and DRAEM [50]. Under the separate case, the baselines' metric is
230 reported in their papers except the metric of US borrowed from [50]. Under the unified case, [US](#),
231 [PSVDD](#), [PaDiM](#), [CutPaste](#), [MKD](#), and [DRAEM](#) are run with the publicly available implementations.

Table 1: **Anomaly detection results with AUROC metric on MVTec-AD [4].** All methods are evaluated under the unified / separate case. In the unified case, the learned model is applied to detect anomalies for all categories *without* fine-tuning.

Category		US [6]	PSVDD [47]	PaDiM [10]	CutPaste [24]	MKD [36]	DRAEM [50]	Ours
Object	Bottle	84.0 / 99.0	85.5 / 98.6	97.9 / 99.9	67.9 / 98.2	98.7 / 99.4	97.5 / 99.2	99.7 \pm 0.04 / 100
	Cable	60.0 / 86.2	64.4 / 90.3	70.9 / 92.7	69.2 / 81.2	78.2 / 89.2	57.8 / 91.8	95.2 \pm 0.84 / 97.6
	Capsule	57.6 / 86.1	61.3 / 76.7	73.4 / 91.3	63.0 / 98.2	68.3 / 80.5	65.3 / 98.5	86.9 \pm 0.73 / 85.3
	Hazelnut	95.8 / 93.1	83.9 / 92.0	85.5 / 92.0	80.9 / 98.3	97.1 / 98.4	93.7 / 100	99.8 \pm 0.10 / 99.9
	Metal Nut	62.7 / 82.0	80.9 / 94.0	88.0 / 98.7	60.0 / 99.9	64.9 / 73.6	72.8 / 98.7	99.2 \pm 0.09 / 99.0
	Pill	56.1 / 87.9	89.4 / 86.1	68.8 / 93.3	71.4 / 94.9	79.7 / 82.7	82.2 / 98.9	93.7 \pm 0.65 / 88.3
	Screw	66.9 / 54.9	80.9 / 81.3	56.9 / 85.8	85.2 / 88.7	75.6 / 83.3	92.0 / 93.9	87.5 \pm 0.57 / 91.9
	Toothbrush	57.8 / 95.3	99.4 / 100	95.3 / 96.1	63.9 / 99.4	75.3 / 92.2	90.6 / 100	94.2 \pm 0.20 / 95.0
	Transistor	61.0 / 81.8	77.5 / 91.5	86.6 / 97.4	57.9 / 96.1	73.4 / 85.6	74.8 / 93.1	99.8 \pm 0.09 / 100
	Zipper	78.6 / 91.9	77.8 / 97.9	79.7 / 90.3	93.5 / 99.9	87.4 / 93.2	98.8 / 100	95.8 \pm 0.51 / 96.7
Texture	Carpet	86.6 / 91.6	63.3 / 92.9	93.8 / 99.8	93.6 / 93.9	69.8 / 79.3	98.0 / 97.0	99.8 \pm 0.02 / 99.9
	Grid	69.2 / 81.0	66.0 / 94.6	73.9 / 96.7	93.2 / 100	83.8 / 78.0	99.3 / 99.9	98.2 \pm 0.26 / 98.5
	Leather	97.2 / 88.2	60.8 / 90.9	99.9 / 100	93.4 / 100	93.6 / 95.1	98.7 / 100	100 \pm 0.00 / 100
	Tile	93.7 / 99.1	88.3 / 97.8	93.3 / 98.1	88.6 / 94.6	89.5 / 91.6	99.8 / 99.6	99.3 \pm 0.14 / 99.0
	Wood	90.6 / 97.7	72.1 / 96.5	98.4 / 99.2	80.4 / 99.1	93.4 / 94.3	99.8 / 99.1	98.6 \pm 0.08 / 97.9
Mean	74.5 / 87.7	76.8 / 92.1	84.2 / 95.5	77.5 / 96.1	81.9 / 87.8	88.1 / 98.0	96.5 \pm 0.08 / 96.6	

Table 2: **Anomaly localization results with AUROC metric on MVTec-AD [4].** All methods are evaluated under the unified / separate case. In the unified case, the learned model is applied to detect anomalies for all categories *without* fine-tuning.

Category		US [6]	PSVDD [47]	PaDiM [10]	FCDD [26]	MKD [36]	DRAEM [50]	Ours
Object	Bottle	67.9 / 97.8	86.7 / 98.1	96.1 / 98.2	56.0 / 97	91.8 / 96.3	87.6 / 99.1	98.1 \pm 0.04 / 98.1
	Cable	78.3 / 91.9	62.2 / 96.8	81.0 / 96.7	64.1 / 90	89.3 / 82.4	71.3 / 94.7	97.3 \pm 0.10 / 96.8
	Capsule	85.5 / 96.8	83.1 / 95.8	96.9 / 98.6	67.6 / 93	88.3 / 95.9	50.5 / 94.3	98.5 \pm 0.01 / 97.9
	Hazelnut	93.7 / 98.2	97.4 / 97.5	96.3 / 98.1	79.3 / 95	91.2 / 94.6	96.9 / 99.7	98.1 \pm 0.10 / 98.8
	Metal Nut	76.6 / 97.2	96.0 / 98.0	84.8 / 97.3	57.5 / 94	64.2 / 86.4	62.2 / 99.5	94.8 \pm 0.09 / 95.7
	Pill	80.3 / 96.5	96.5 / 95.1	87.7 / 95.7	65.9 / 81	69.7 / 89.6	94.4 / 97.6	95.0 \pm 0.16 / 95.1
	Screw	90.8 / 97.4	74.3 / 95.7	94.1 / 98.4	67.2 / 86	92.1 / 96.0	95.5 / 97.6	98.3 \pm 0.08 / 97.4
	Toothbrush	86.9 / 97.9	98.0 / 98.1	95.6 / 98.8	60.8 / 94	88.9 / 96.1	97.7 / 98.1	98.4 \pm 0.03 / 97.8
	Transistor	68.3 / 73.7	78.5 / 97.0	92.3 / 97.6	54.2 / 88	71.7 / 76.5	64.5 / 90.9	97.9 \pm 0.19 / 98.7
	Zipper	84.2 / 95.6	95.1 / 95.1	94.8 / 98.4	63.0 / 92	86.1 / 93.9	98.3 / 98.8	96.8 \pm 0.24 / 96.0
Texture	Carpet	88.7 / 93.5	78.6 / 92.6	97.6 / 99.0	68.6 / 96	95.5 / 95.6	98.6 / 95.5	98.5 \pm 0.01 / 98.0
	Grid	64.5 / 89.9	70.8 / 96.2	71.0 / 97.1	65.8 / 91	82.3 / 91.8	98.7 / 99.7	96.5 \pm 0.04 / 94.6
	Leather	95.4 / 97.8	93.5 / 97.4	84.8 / 99.0	66.3 / 98	96.7 / 98.1	97.3 / 98.6	98.8 \pm 0.03 / 98.3
	Tile	82.7 / 92.5	92.1 / 91.4	80.5 / 94.1	59.3 / 91	85.3 / 82.8	98.0 / 99.2	91.8 \pm 0.10 / 91.8
	Wood	83.3 / 92.1	80.7 / 90.8	89.1 / 94.1	53.3 / 88	80.5 / 84.8	96.0 / 96.4	93.2 \pm 0.08 / 93.4
Mean	81.8 / 93.9	85.6 / 95.7	89.5 / 97.4	63.3 / 92	84.9 / 90.7	87.2 / 97.3	96.8 \pm 0.02 / 96.6	

232 **Quantitative results of anomaly detection on MVTec-AD [4]** are shown in Tab. 1. Though
 233 all baselines achieve excellent performances under the separate case, their performances drop
 234 dramatically under the unified case. The previous SOTA, DRAEM, a reconstruction-based method
 235 trained by pseudo-anomaly, suffers from a drop of near 10%. For another strong baseline, CutPaste,
 236 a pseudo-anomaly approach, the drop is as large as 18.6%. However, our UniAD has almost no
 237 performance drop from the separate case (96.6%) to the unified case (96.5%). Moreover, we beat the
 238 best competitor, DRAEM, by a dramatically large margin (8.4%), demonstrating our superiority.

239 4.3 Anomaly localization on MVTec-AD

240 **Setup and baselines.** Anomaly localization aims to localize anomalous regions in an anomalous
 241 image. MVTec-AD [4] is chosen as the benchmark dataset. The setup is the same as that in Sec. 4.2.
 242 Besides the competitors in Sec. 4.2, FCDD [26] is included, whose metric under the separate case is
 243 reported in its paper. Under the unified case, we run FCDD with the implementation: **FCDD**.

244 **Quantitative results of anomaly localization on MVTec-AD [4]** are reported in Tab. 2. Similar
 245 to Sec. 4.2, switching from the separate case to the unified case, the performance of all competitors

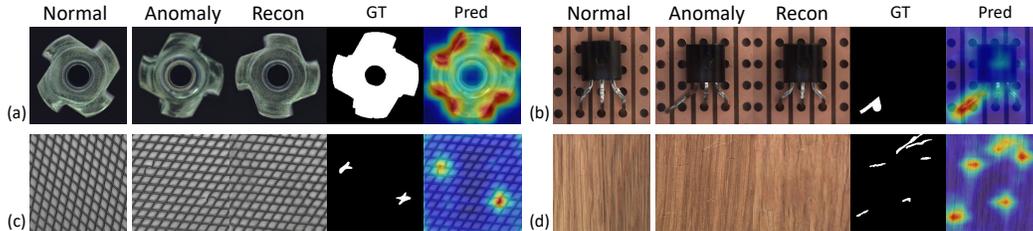


Figure 5: **Qualitative results** for anomaly localization on MVTec-AD [4]. From left to right: normal sample as the reference, anomaly, our reconstruction, ground-truth, and our predicted anomaly map. The approach to visualizing reconstruction is the same as the one used in Fig. 2.

Table 3: **Anomaly detection results with AUROC metric on CIFAR-10** [22] under the unified case. Here, {01234} means samples from class 0, 1, 2, 3, 4 are borrowed as the normal ones.

Normal Indices	US [6]	FCDD [26]	FCDD+OE [26]	PANDA [32]	MKD [36]	Ours
{01234}	51.3	55.0	71.8	66.6	64.2	80.1 \pm 0.08
{56789}	51.3	50.3	73.7	73.2	69.3	73.8 \pm 0.08
{02468}	63.9	59.2	85.3	77.1	76.4	88.8 \pm 0.20
{13579}	56.8	58.5	85.0	72.9	78.7	85.6 \pm 0.10
Mean	55.9	55.8	78.9	72.4	72.1	82.1 \pm 0.08

246 drops significantly. For example, the performance of US, an important distillation-based baseline,
 247 decreases by 12.1%. FCDD, a pseudo-anomaly approach, suffers from a dramatic drop of 28.7%,
 248 reflecting the pseudo-anomaly is not suitable for the unified case. However, our UniAD even gains a
 249 slight improvement from the separate case (96.6%) to the unified case (96.8%), proving the suitability
 250 of our UniAD for the unified case. Moreover, we significantly surpass the strongest baseline, PaDiM,
 251 by 7.3%. This significant improvement reflects the effectiveness of our model.

252 **Qualitative results for anomaly localization on MVTec-AD** [4] are illustrated in Fig. 5. For both
 253 global (Fig. 5a) and local (Fig. 5b) structural anomalies, both scattered texture perturbations (Fig. 5c)
 254 and multiple texture scratches (Fig. 5d), our method could successfully reconstruct anomalies to their
 255 corresponding normal samples, then accurately localize anomalous regions through reconstruction
 256 differences. More qualitative results are given in *Supplementary Material*.

257 4.4 Anomaly detection on CIFAR-10

258 **Setup.** To further verify the effectiveness of our UniAD, we extend CIFAR-10 [22] to the unified
 259 case, which consists of four combinations. For each combination, five categories together serve
 260 as normal samples, while other categories are viewed as anomalies. The class indices of the four
 261 combinations are {01234}, {56789}, {02468}, {13579}. Here, {01234} means the normal samples
 262 include images from class 0, 1, 2, 3, 4, and similar for others. Note that the class index is obtained by
 263 sorting the class names of 10 classes. The setup of the model is detailed in *Supplementary Material*.

264 **Baselines.** US [6], FCDD [26], FCDD+OE [26], PANDA [32], and MKD [36] serve as competitors.
 265 US, FCDD, FCDD+OE, PANDA, and MKD are run with the publicly available implementations.

266 **Quantitative results of anomaly detection on CIFAR-10** [22] are shown in Tab. 3. When five
 267 classes together serve as normal samples, two recent baselines, US and FCDD, almost lose their
 268 ability to detect anomalies. When utilizing 10000 images sampled from CIFAR-100 [22] as auxiliary
 269 Outlier Exposure (OE), FCDD+OE improves the performance by a large margin. We still stably
 270 outperform FCDD+OE by 3.2% without the help of OE, indicating the efficacy of our UniAD.

271 4.5 Ablation studies

272 To verify the effectiveness of the proposed modules and the selection of hyperparameters, we
 273 implement extensive ablation studies on MVTec-AD [4] under the unified case.

274 **Layer-wise query.** Tab. 4a verifies our assertion that the query embedding is of vital significance.
 275 1) Without query embedding, meaning the encoder embeddings are directly input to the decoder,
 276 the performance is the worst. 2) Adding only one query embedding to the first decoder layer (*i.e.*,
 277 vanilla transformer [41]) promotes the performance dramatically by 13.4% and 18.1% in anomaly

Table 4: **Ablation studies with AUROC metric on MVTec-AD [4].** Default settings are in blue.

(a) Layer-wise query, NMA, & FJ						(b) Layer Number of Encoder & Decoder					
w/o q.	1 q.	Layer-wise q.	NMA	FJ	Loc.	Det.	#Enc, #Dec	Vanilla [41]		Ours	
								Loc.	Det.	Loc.	Det.
✓	-	-	-	-	79.4	69.5					
-	✓	-	-	-	92.8	87.6	4, 0	79.2	69.8	96.0	94.9
-	-	✓	-	-	96.5	95.0	0, 4	88.3	80.5	96.3	96.1
-	✓	-	✓	-	96.3	96.1	2, 2	90.6	84.7	96.0	95.1
-	✓	-	-	✓	95.8	95.0	4, 4	92.8	87.6	96.8	96.5
-	-	✓	✓	✓	96.8	96.5	6, 6	91.9	86.1	96.7	96.5

(c) Neighbor Size in NMA			(d) Where to Add NMA			(e) Jitter Scale α in FJ			(f) Jitter Prob. p in FJ		
Size	Loc.	Det.	Place	Loc.	Det.	α	Loc.	Det.	p	Loc.	Det.
1×1	96.3	94.6	Enc	96.3	95.8	5	96.7	96.1	0.25	96.5	95.6
5×5	96.8	96.4	Enc+Dec1	96.8	96.4	10	96.7	96.4	0.50	96.7	95.8
7×7	96.8	96.5	Enc+Dec2	96.7	96.5	20	96.8	96.5	0.75	96.7	96.3
9×9	96.7	96.3	All	96.8	96.5	30	96.6	95.7	1	96.8	96.5

278 localization and detection, respectively. 3) With layer-wise query embedding in each decoder layer,
 279 pixel-level and image-level AUROC is further improved by 3.7% and 7.4%, respectively.

280 **Layer number.** We conduct experiments to investigate the influence of layer number, as shown
 281 in Tab. 4b. 1) No matter with which combination, our model outperforms vanilla transformer by a
 282 large margin, reflecting the effectiveness of our design. 2) The best performance is achieved with a
 283 moderate layer number: 4Enc+4Dec. A larger layer number like 6Enc+6Dec does not bring further
 284 promotion, which may be because more layers are harder to train.

285 **Neighbor masked attention.** 1) The effectiveness of NMA is proven in Tab. 4a. Under the case
 286 of one query embedding, adding NMA brings promotion by 3.5% for localization and 8.5% for
 287 detection. 2) The neighbor size of NMA is selected in Tab. 4c. 1×1 neighbor size is the worst,
 288 because 1×1 is too small to prevent the information leak, thus the recovery could be completed by
 289 copying neighbor regions. A larger neighbor size ($\geq 5 \times 5$) is obviously much better, and the best
 290 one is selected as 7×7. 3) We also study the place to add NMA in Tab. 4d. Only adding NMA in the
 291 encoder (Enc) is not enough. The performance could be stably improved when further adding NMA
 292 in the first or second attention in the decoder (Enc+Dec1, Enc+Dec2) or both (All). This reflects that
 293 the full attention of the decoder also contributes to the information leak.

294 **Feature jittering.** 1) Tab. 4a confirms the efficacy of FJ. With one query embedding as the baseline,
 295 introducing FJ could bring an increase of 3.0% for localization and 7.4% for detection, respectively.
 296 2) According to Tab. 4e, the jittering scale, α , is chosen as 20. A larger α (*i.e.*, 30) disturbs the feature
 297 too much, degrading the results. 3) In Tab. 4f, the jittering probability, p , is studied. In essence, the
 298 task would be a denoising task with feature jittering, and be a reconstruction task without feature
 299 jittering. The results show that the full denoising task (*i.e.*, $p = 1$) is the best.

300 5 Conclusion

301 In this work, we propose UniAD that unifies anomaly detection regarding multiple classes. For such a
 302 challenging task, we assist the model against learning an “identical shortcut” with three improvements.
 303 First, we confirm the effectiveness of the learnable query embedding and carefully tailor a layer-wise
 304 query decoder to help model the complex distribution of multi-class data. Second, we come up with a
 305 neighbor masked attention module to avoid the information leak from the input to the output. Third,
 306 we propose feature jittering that helps the model less sensitive to the input perturbations. Under the
 307 unified task setting, our method achieves state-of-the-art performance on MVTec-AD and CIFAR-10
 308 datasets, significantly outperforming existing alternatives.

309 **Discussion.** In this work, different kinds of objects are handled without being distinguished. We have
 310 not used the category labels that may help the model better fit multi-class data. How to incorporate
 311 the unified model with category labels should be further studied. In practical uses, normal samples are
 312 not as consistent as those in MVTec-AD, often manifest themselves in some diversity. Our UniAD
 313 could handle all 15 categories in MVTec-AD, hence would be more suitable for real scenes. **However,**
 314 **anomaly detection may be used for video surveillance, which may infringe personal privacy.**

References

- 315 [1] F. Ahmed and A. Courville. Detecting semantic anomalies. In *Assoc. Adv. Artif. Intell.*, 2020.
- 316 [2] S. Akcay, A. Atapour-Abarghouei, and T. P. Breckon. GANomaly: Semi-supervised anomaly detection via
317 adversarial training. In *Asian Conf. Comput. Vis.*, 2018.
- 318 [3] Y. Bengio, L. Yao, G. Alain, and P. Vincent. Generalized denoising auto-encoders as generative models. In
319 *Adv. Neural Inform. Process. Syst.*, 2013.
- 320 [4] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger. MVTEC AD—A comprehensive real-world dataset for
321 unsupervised anomaly detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.
- 322 [5] P. Bergmann, S. Löwe, M. Fauser, D. Sattlegger, and C. Steger. Improving unsupervised defect segmentation
323 by applying structural similarity to autoencoders. In *International Joint Conference on Computer Vision,*
324 *Imaging and Computer Graphics Theory and Applications (VISIGRAPP)*, 2019.
- 325 [6] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger. Uninformed students: Student-teacher anomaly
326 detection with discriminative latent embeddings. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.
- 327 [7] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko. End-to-end object detection
328 with transformers. In *Eur. Conf. Comput. Vis.*, 2020.
- 329 [8] A.-S. Collin and C. De Vleeschouwer. Improved anomaly detection by training an autoencoder with skip
330 connections on images corrupted with stain-shaped noise. In *Int. Conf. Pattern Recog.*, 2021.
- 331 [9] L. Deecke, L. Ruff, R. A. Vandermeulen, and H. Bilen. Transfer-based semantic anomaly detection. In *Int.*
332 *Conf. Mach. Learn.*, 2021.
- 333 [10] T. Defard, A. Setkov, A. Loesch, and R. Audigier. PaDim: A patch distribution modeling framework for
334 anomaly detection and localization. In *Int. Conf. Pattern Recog.*, 2021.
- 335 [11] D. Dehaene, O. Frigo, S. Combrexelle, and P. Eline. Iterative energy-based projection on a normal data
336 manifold for anomaly localization. In *Int. Conf. Learn. Represent.*, 2019.
- 337 [12] H. Deng and X. Li. Anomaly detection via reverse distillation from one-class embedding. In *IEEE Conf.*
338 *Comput. Vis. Pattern Recog.*, 2022.
- 339 [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical image
340 database. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2009.
- 341 [14] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers
342 for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 343 [15] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,
344 M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image
345 recognition at scale. In *Int. Conf. Learn. Represent.*, 2021.
- 346 [16] T. Fernando, H. Gammulle, S. Denman, S. Sridharan, and C. Fookes. Deep learning for medical anomaly
347 detection—a survey. *ACM Computing Surveys (CSUR)*, 2021.
- 348 [17] I. Golan and R. El-Yaniv. Deep anomaly detection using geometric transformations. In *Adv. Neural Inform.*
349 *Process. Syst.*, 2018.
- 350 [18] D. Gong, L. Liu, V. Le, B. Saha, M. R. Mansour, S. Venkatesh, and A. v. d. Hengel. Memorizing normality
351 to detect anomaly: Memory-augmented deep autoencoder for unsupervised anomaly detection. In *Int.*
352 *Conf. Comput. Vis.*, 2019.
- 353 [19] J. Hou, Y. Zhang, Q. Zhong, D. Xie, S. Pu, and H. Zhou. Divide-and-assemble: Learning block-wise
354 memory for unsupervised anomaly detection. In *Int. Conf. Comput. Vis.*, 2021.
- 355 [20] L. Ilya and H. Frank. Decoupled weight decay regularization. In *Int. Conf. Learn. Represent.*, 2019.
- 356 [21] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- 357 [22] A. Krizhevsky. Learning multiple layers of features from tiny images. *Master’s thesis, University of Tront*,
358 2009.
- 359 [23] G. Kwon, M. Prabhushankar, D. Temel, and G. AlRegib. Backpropagated gradient representations for
360 anomaly detection. In *Eur. Conf. Comput. Vis.*, 2020.
- 361

- 362 [24] C.-L. Li, K. Sohn, J. Yoon, and T. Pfister. CutPaste: Self-supervised learning for anomaly detection and
363 localization. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.
- 364 [25] W. Liu, R. Li, M. Zheng, S. Karanam, Z. Wu, B. Bhanu, R. J. Radke, and O. Camps. Towards visually
365 explaining variational autoencoders. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.
- 366 [26] P. Liznerski, L. Ruff, R. A. Vandermeulen, B. J. Franks, M. Kloft, and K.-R. Müller. Explainable deep
367 one-class classification. In *Int. Conf. Learn. Represent.*, 2021.
- 368 [27] P. Mishra, R. Verk, D. Fornasier, C. Piciarelli, and G. L. Foresti. VT-ADL: A vision transformer network
369 for image anomaly detection and localization. In *International Symposium on Industrial Electronics*, 2021.
- 370 [28] H. Park, J. Noh, and B. Ham. Learning memory-guided normality for anomaly detection. In *IEEE Conf.*
371 *Comput. Vis. Pattern Recog.*, 2020.
- 372 [29] P. Perera, R. Nallapati, and B. Xiang. OCGAN: One-class novelty detection using GANs with constrained
373 latent representations. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.
- 374 [30] J. Pirnay and K. Chai. Inpainting transformer for anomaly detection. *arXiv preprint arXiv:2104.13897*,
375 2021.
- 376 [31] M. Pourreza, B. Mohammadi, M. Khaki, S. Bouindour, H. Snoussi, and M. Sabokrou. G2D: generate to
377 detect anomaly. In *IEEE Winter Conf. Appl. Comput. Vis.*, 2021.
- 378 [32] T. Reiss, N. Cohen, L. Bergman, and Y. Hoshen. Panda: Adapting pretrained features for anomaly detection
379 and segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.
- 380 [33] O. Rippel, P. Mertens, and D. Merhof. Modeling the distribution of normal data in pretrained deep features
381 for anomaly detection. In *Int. Conf. Pattern Recog.*, 2021.
- 382 [34] L. Ruff, R. Vandermeulen, N. Goernitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, and M. Kloft.
383 Deep one-class classification. In *Int. Conf. Mach. Learn.*, 2018.
- 384 [35] M. Sabokrou, M. Khalooei, M. Fathy, and E. Adeli. Adversarially learned one-class classifier for novelty
385 detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018.
- 386 [36] M. Salehi, N. Sadjadi, S. Baselizadeh, M. H. Rohban, and H. R. Rabiee. Multiresolution knowledge
387 distillation for anomaly detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.
- 388 [37] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson. Estimating the support of a
389 high-dimensional distribution. *Neural Computation*, 2001.
- 390 [38] Y. Shi, J. Yang, and Z. Qi. Unsupervised anomaly segmentation via deep feature reconstruction.
391 *Neurocomputing*, 2021.
- 392 [39] M. Tan and Q. Le. EfficientNet: Rethinking model scaling for convolutional neural networks. In *Int. Conf.*
393 *Mach. Learn.*, 2019.
- 394 [40] D. M. Tax and R. P. Duin. Support vector data description. *Machine Learning*, 2004.
- 395 [41] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin.
396 Attention is all you need. In *Adv. Neural Inform. Process. Syst.*, 2017.
- 397 [42] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol. Extracting and composing robust features with
398 denoising autoencoders. In *Int. Conf. Mach. Learn.*, 2008.
- 399 [43] G. Wang, S. Han, E. Ding, and D. Huang. Student-teacher feature pyramid matching for anomaly detection.
400 *Brit. Mach. Vis. Conf.*, 2021.
- 401 [44] S. Wang, L. Wu, L. Cui, and Y. Shen. Glancing at the patch: Anomaly localization with global and local
402 feature comparison. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.
- 403 [45] Y. Xia, Y. Zhang, F. Liu, W. Shen, and A. L. Yuille. Synthesize then compare: Detecting failures and
404 anomalies for semantic segmentation. In *Eur. Conf. Comput. Vis.*, 2020.
- 405 [46] X. Yan, H. Zhang, X. Xu, X. Hu, and P.-A. Heng. Learning semantic context from normal samples for
406 unsupervised anomaly detection. In *Assoc. Adv. Artif. Intell.*, 2021.
- 407 [47] J. Yi and S. Yoon. Patch SVDD: Patch-level SVDD for anomaly detection and segmentation. In *Asian*
408 *Conf. Comput. Vis.*, 2020.

- 409 [48] L. Yunseung and K. Pilsung. AnoViT: Unsupervised anomaly detection and localization with vision
410 transformer-based encoder-decoder. *arXiv preprint arXiv:2203.10808*, 2022.
- 411 [49] M. Z. Zaheer, J.-h. Lee, M. Astrid, and S.-I. Lee. Old is gold: Redefining the adversarially learned one-class
412 classifier training paradigm. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.
- 413 [50] V. Zavrtanik, M. Kristan, and D. Skočaj. DRAEM-A discriminatively trained reconstruction embedding
414 for surface anomaly detection. In *Int. Conf. Comput. Vis.*, 2021.
- 415 [51] K. Zhou, Y. Xiao, J. Yang, J. Cheng, W. Liu, W. Luo, Z. Gu, J. Liu, and S. Gao. Encoding structure-texture
416 relation with P-Net for anomaly detection in retinal images. In *Eur. Conf. Comput. Vis.*, 2020.

417 Checklist

- 418 1. For all authors...
- 419 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
420 contributions and scope? [Yes]
- 421 (b) Did you describe the limitations of your work? [Yes] See Sec. 5.
- 422 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 423 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
424 them? [Yes]
- 425 2. If you are including theoretical results...
- 426 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 427 (b) Did you include complete proofs of all theoretical results? [N/A]
- 428 3. If you ran experiments...
- 429 (a) Did you include the code, data, and instructions needed to reproduce the main
430 experimental results (either in the supplemental material or as a URL)? [Yes] See
431 Sec. 4 and the supplemental material.
- 432 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
433 were chosen)? [Yes] See Sec. 4.1 & Sec. 4.4 for data splits, and Sec. 4.5 for the choice
434 of hyperparameters.
- 435 (c) Did you report error bars (e.g., with respect to the random seed after running
436 experiments multiple times)? [Yes] See Tab. 1, Tab. 2, and Tab. 3.
- 437 (d) Did you include the total amount of compute and the type of resources used (e.g., type
438 of GPUs, internal cluster, or cloud provider)? [Yes] See **Setup** in Sec. 4.2.
- 439 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 440 (a) If your work uses existing assets, did you cite the creators? [Yes] See Sec. 4.
- 441 (b) Did you mention the license of the assets? [N/A]
- 442 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 443
- 444 (d) Did you discuss whether and how consent was obtained from people whose data you’re
445 using/curating? [N/A]
- 446 (e) Did you discuss whether the data you are using/curating contains personally identifiable
447 information or offensive content? [N/A]
- 448 5. If you used crowdsourcing or conducted research with human subjects...
- 449 (a) Did you include the full text of instructions given to participants and screenshots, if
450 applicable? [N/A]
- 451 (b) Did you describe any potential participant risks, with links to Institutional Review
452 Board (IRB) approvals, if applicable? [N/A]
- 453 (c) Did you include the estimated hourly wage paid to participants and the total amount
454 spent on participant compensation? [N/A]