

Prototypical Residual Networks for Anomaly Detection and Localization

Hui Zhang^{1,2} Zuxuan Wu^{1,2} Zheng Wang³ Zhineng Chen^{1,2*} Yu-Gang Jiang^{1,2}

¹Shanghai Key Lab of Intell. Info. Processing, School of CS, Fudan University

²Shanghai Collaborative Innovation Center of Intelligent Visual Computing

³School of Computer Science, Zhejiang University of Technology

Abstract

Anomaly detection and localization are widely used in industrial manufacturing for its efficiency and effectiveness. Anomalies are rare and hard to collect and supervised models easily over-fit to these seen anomalies with a handful of abnormal samples, producing unsatisfactory performance. On the other hand, anomalies are typically subtle, hard to discern, and of various appearance, making it difficult to detect anomalies and let alone locate anomalous regions. To address these issues, we propose a framework called Prototypical Residual Network (PRN), which learns feature residuals of varying scales and sizes between anomalous and normal patterns to accurately reconstruct the segmentation maps of anomalous regions. PRN mainly consists of two parts: multi-scale prototypes that explicitly represent the residual features of anomalies to normal patterns; a multi-size self-attention mechanism that enables variable-sized anomalous feature learning. Besides, we present a variety of anomaly generation strategies that consider both seen and unseen appearance variance to enlarge and diversify anomalies. Extensive experiments on the challenging and widely used MVTec AD benchmark show that PRN outperforms current state-of-the-art unsupervised and supervised methods. We further report SOTA results on three additional datasets to demonstrate the effectiveness and generalizability of PRN.

1. Introduction

The human cognition and visual system has an inherent ability to perceive anomalies [53]. Not only can humans distinguish between defective and non-defective images, but they can also point to the location of anomalies even if they have seen none or only a limited number of anomalies. Anomaly detection (image-level binary classification) and anomaly localization (pixel-level binary classification) are introduced for the same purpose, and have been widely used

* Corresponding author.

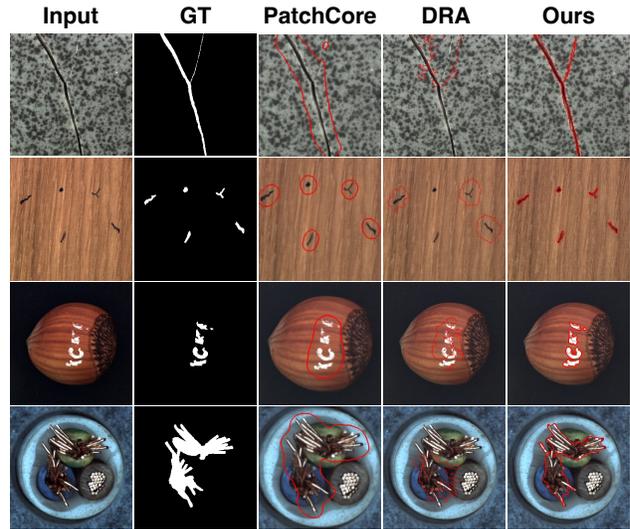


Figure 1. Anomaly detection and localization examples on MVTec [4]. Compared with the unsupervised method PatchCore [41] and the supervised method DRA [13], the proposed PRN is able to locate the anomalous regions more accurately.

in various scenarios due to their efficiency and remarkable accuracy, including industrial defect detection [4, 7, 34, 61], medical image analysis [52] and video surveillance [32].

Given its importance, a significant amount of work has been devoted to anomaly detection and anomaly localization, but few have addressed both detection and localization problems well at the same time. We argue that real-world anomalous data weaken these models mainly in three aspects: I) the amount of abnormal samples is limited and significant fewer than normal samples, producing data distributions that lead to a naturally **imbalanced learning** problem; II) anomalies are typically subtle and hard to discern, since normal patterns still dominate the anomalous image; **identifying abnormal regions** out of the whole image is the key to anomaly detection and localization; III) the appearance of anomalies varies significantly, *i.e.*, abnormal regions can take on a variety of sizes, shapes and numbers, and such **appearance variations**

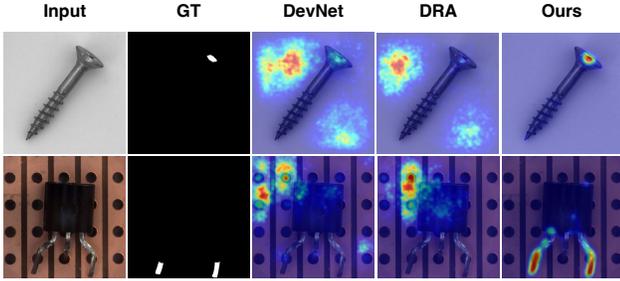


Figure 2. Indecipherable problem of supervised methods DevNet [35] and DRA [13]. Both images are detected as anomalous. Other methods mistakenly highlight normal regions rather than defect regions, whereas PRN correctly pinpoints the defect regions.

make it challenging to well-localizing all the anomalies.

Without adequate anomalies for training, unsupervised models become the de facto dominant approaches, which get rid of the imbalance problem by learning the distribution of normal samples [5, 9, 10, 12, 18, 25, 41–43, 47, 67] or generating sufficient synthetic anomalies [26, 28, 50, 63, 68]. However, these methods are opaque to genuine anomalies, resulting in implicit decisions that may induce many false negatives and false positives. Besides, unsupervised methods rely heavily on the quality of normal samples, and thus are not robust enough and perform poorly on uncalibrated or noisy datasets [19]. As shown in Fig. 1, unsupervised models predict broad regions around the anomaly. We attribute this problem to less discriminative abilities of these methods.

Recently, several supervised methods [13, 35, 46] are introduced. DeepSAD [46] enlarges the margin between the anomaly and the one-class center in the latent space to obtain more compact one-class descriptors by limit seen anomalies. DRA [13] and DevNet [35] formulate anomaly detection as a multi-instance learning (MIL) problem, scoring an image as anomaly if any image patch is a defect region. MIL-based methods enforce the learning at fine-grained image patch level, which effectively reduces the interference of normal patches in the anomalous images. Yet, these approaches typically struggle to accurately locate all anomalous regions with image-level supervision, as shown in Fig. 1. In particular, when the anomalous regions only occupy a tiny part of image patches, image-level representation may be dominated by the normal regions and disregards tiny anomalous, which may cause inconsistent image-level and pixel level performance as shown in Table 1. Furthermore, as shown in Fig. 2, these methods also encounter uninterpretable problems when making decisions.

In this paper, we propose a framework called Prototypical Residual Network (PRN) as an effective remedy for afore-said issues on anomaly detection and localization. First, we propose multi-scale prototypes to represent normal patterns. In contrast to previous methods for constructing normal pat-

terns from concatenated feature memory [41] or random sampled feature maps [63], we construct normal patterns with prototypes of intermediate feature maps of different scales, thereby preserving the spatial information and providing precise and representative normal patterns. Further, we obtain the feature map residuals via the deviation between the anomalous image and the closest prototype at each scale, and we add multi-scale fusion blocks to exchange information across different scales. Second, since the appearance of anomaly regions varies a lot, it is necessary to learn relationships among patches from multiple receptive fields. Thus, we introduce a multi-size self-attention [33, 55, 58, 59] mechanism, which operates on patches of different receptive fields to detect patch-level inconsistencies at different sizes. Finally, unlike previous methods [13, 35] that use image-level supervision for training, our model learns to reconstruct the anomaly segmentation map with pixel-level supervision, which focuses more on the anomalous regions and preserves better generalization. Besides, we put forward a variety of anomaly generation strategies that efficiently mitigate the impact of data imbalance and enrich the anomaly appearance. With the proposed modules, our method achieves more accurate localization than previous unsupervised and supervised methods, as shown in Fig. 1 and Fig. 2.

The main contributions of this paper are summarized as follows:

- We propose a novel Prototypical Residual Networks for anomaly detection and localization. Equipped with multi-scale prototypes and the multi-size self-attention mechanism, PRN learns residual representations among multi-scale feature maps and within the multi-size receptive fields at each scale.
- We present a variety of anomaly generation strategies that considering both seen and unseen appearance variance to enlarge and diversify anomalies.
- We perform extensive experiments on four datasets to show that our approach achieves new SOTA anomaly detection performance and outperforms current SOTA in anomaly localization performance by a large margin.

2. Related Work

Unsupervised Approaches. Unsupervised paradigm assumes that only normal data is available during training [36, 44, 53]. Auto-Encoder based methods [3, 6, 15, 65] rely on the hypothesis that the model is trained to reconstruct normal regions well but fails for abnormal regions. Although localization results based on the difference between the input and the reconstructed image are often intuitive and interpretable, their performance is limited. Generative models are introduced to obtain better reconstruction performance. However, the generation effect of VAE [11, 11, 30]

or GAN [2, 22, 27, 48, 49, 66] over normal areas in the image is poor, leading to coarse reconstruction and false detection. Normalizing flows based methods [18, 42, 43, 67] learn bijective transformations between data distributions and well-defined densities, however the computational cost of these approaches is significant. Knowledge distillation-based methods [5, 12, 47, 56] transform the anomaly detection task into a feature comparison between teacher and student networks. Deep feature modeling-based methods [9, 10, 23, 25, 34, 41, 70] build a feature space for input images and then detect and localize anomalies by comparing the features. Self-supervised learning-based methods [38, 39, 64] designed proxy tasks such as predicting or recovering hidden regions or properties in input images [53]. One-class classification based methods [45, 51, 64] aim to map training images or patches to a small hypersphere in the feature space. These approaches address the imbalance problem by being opaque to anomalous samples, but suffer from implicit decisions that result in subpar performance on subtle and challenging anomalies.

Supervised Approaches. A recent emerging trend focuses on supervised anomaly detection by leveraging seen anomalies to increase the differentiation between anomalous and normal samples. Some existing works [16, 31, 46] are learned with a minority of anomaly based on one-class classification metric. Some anomaly-focused deviation losses proposed in [35, 69] mitigate the bias derived from the seen anomalies. A multi-head model is introduced in [13] to learn disentangled anomaly representations, where each head is dedicated to capturing a specific type of anomaly. Due to the imbalanced learning problem, these methods are prone to over-fitting to seen anomalies and fail to generalize to unseen anomalies, resulting in poor anomaly detection performance. In addition, image-level representations may be dominated by normal regions while disregarding the representations of subtle anomaly regions, resulting in an inability to accurately localize anomalies that come in a variety of sizes, shapes and numbers.

3. Method

Together with the proposed anomaly generation strategies (Appendix A.1) that retain a balanced data distribution, we propose the Prototypical Residual Network (PRN) to reconstruct the segmentation map for anomaly detection and localization. Overall, we adopt a U-Net [40]-like architecture as shown in Fig. 3. The encoder is a pre-trained ResNet-18 [21], and the decoder consists of upsampling and convolution blocks. The skip-connection branches of PRN are equipped with the proposed Multi-scale Prototypes (MP, Sec. 3.1), Multi-scale Fusion blocks (MF, Sec. 3.2) and a Multi-size Self-Attention mechanism (MSA, Sec. 3.3). In the following, we will concretely describe each part.

3.1. Multi-scale Prototypes

Prototype Initialization. We define \mathcal{X}_N to be the set of all normal samples during training ($\forall x \in \mathcal{X}_N : y_x = 0$). y_x denotes that if an image x is normal (0) or abnormal (1). Following [10, 41], we use a network pre-trained on ImageNet to obtain feature maps of the input image at different scales. We use $\mathcal{F}_{i,j} = \mathcal{F}_j(x_i)$ ($j \in \{1, 2, 3, 4\}$) to denote the j -th block output of input x_i from a ResNet-like architecture such as ResNet-18 [21]. Assume the feature map $\mathcal{F}_{i,j} \in \mathbb{R}^{c^j \times h^j \times w^j}$ to be a tensor of depth c^j , height h^j and width w^j . Firstly, the j -th scale prototypes $\mathcal{P}_j \in \mathbb{R}^{K \times c^j \times h^j \times w^j}$ are K feature maps randomly sampled from $\mathcal{F}_j(\mathcal{X}_N)$, and are updated by k-means clustering [20]. L2 distance is used to calculate the distance between two feature maps. As the number of normal samples in different datasets varies considerably, to have a suitable amount of prototypes, we set the number of prototypes to a certain ratio of the number of normal samples. As a result, the value of K varies by datasets. The ablation on the proportion number is detailed in Sec. 4.3, and is typically 10%. Three scales of prototypes are employed ($j \in \{1, 2, 3\}$). Model parameters are frozen during clustering. After clustering, the prototypes $\mathcal{P}_j \in \mathbb{R}^{K \times c^j \times h^j \times w^j}$ at each scale remain unchanged during subsequent model training.

Residual Representation. Given the i -th input image and its corresponding feature map $\mathcal{F}_{i,j}$ at j -th block, we can find the closest prototype \mathcal{P}_j^* at j -th scale by calculating the L2 distance between $\mathcal{F}_{i,j}$ and each of the prototypes \mathcal{P}_j . We define the anomalous residual representation of $\mathcal{F}_{i,j}$ to its closest prototype as

$$\begin{aligned} \mathcal{D}_{i,j} &= D(\mathcal{F}_{i,j} - \mathcal{P}_j^*), \\ \text{s.t. } \mathcal{P}_j^* &= \arg \min_{\mathcal{P}_j^k \in \mathcal{P}_j} \|\mathcal{F}_{i,j} - \mathcal{P}_j^k\|_2 \end{aligned} \quad (1)$$

where $D(\cdot, \cdot)$ implements the element-wise Euclidean distance between two tensors, $\mathcal{D}_{i,j} \in \mathbb{R}^{c^j \times h^j \times w^j}$ is the residual from the nearest cluster prototype \mathcal{P}_j^* . Note that the input sample can match distinct prototypes at different scales, as the prototypes are learned independently at each scale.

3.2. Multi-scale Fusion

To enable information exchanging across multi-scale representations, we propose to use Multi-scale Fusion blocks (MF) inspired by [17, 57]. As shown in Fig. 4, the fused output feature map is the sum of the transformed representations of three input feature maps. The feature map $\mathcal{F}_{i,j}^*$ is fused with others as follows:

$$\mathcal{F}_{i,j}^* = f_{1j}(\mathcal{F}_{i,1}) + f_{2j}(\mathcal{F}_{i,2}) + f_{3j}(\mathcal{F}_{i,3}) \quad (2)$$

The choice of the transform function $f_{rj}(\cdot)$ depends on the input feature map index r and the output feature map index

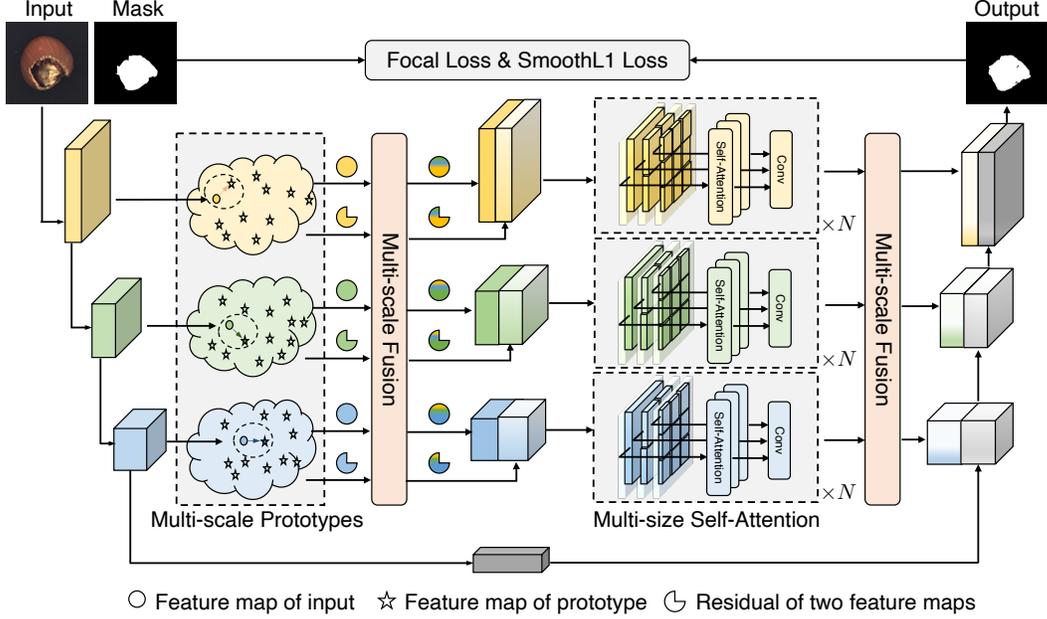


Figure 3. An overview of the proposed Prototypical Residual Network. Anomalous feature residuals of inputs are obtained via Multi-scale Prototypes for each scale feature map from the nearest cluster prototype. Feature maps and residuals at different scales are separately fused by Multi-scale Fusion blocks. Multi-size Self-Attention learns feature residuals on patches of different sizes at each scale, which are further enhanced by another Multi-scale Fusion block. Please see text for details.

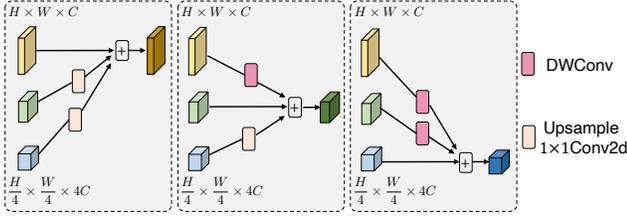


Figure 4. Three feature maps of different scales are fused by a multi-scale fusion block.

j ($r, j \in \{1, 2, 3\}$). If $r = j$, $f_{rj}(\mathcal{F}_{i,r}) = \mathcal{F}_{i,r}$. If $r < j$, $f_{rj}(\mathcal{F}_{i,r})$ down-samples the input feature map $\mathcal{F}_{i,r}$ through depth-wise separable convolutions with a stride of 2^{j-r} , a kernel size of $2^{j-r} + 1$ and a padding of 2^{j-r-1} . If $r > j$, $f_{rj}(\mathcal{F}_{i,r})$ up-samples the input feature map $\mathcal{F}_{i,r}$ through a bilinear up-sampling followed by a 1×1 convolution. The anomalous residual representation $\mathcal{D}_{i,j}$ also follows the fusion paradigm in Fig. 4 and is concatenated with $\mathcal{F}_{i,j}^*$ along the depth dimension to obtain $\mathcal{C}_{i,j}^* \in \mathbb{R}^{2c^j \times h^j \times w^j}$.

3.3. Multi-size Self-Attention

As the anomalous regions vary in magnitude, to further detect local inconsistencies in the concatenated feature maps $\mathcal{C}_{i,j}^*$, we introduce a Multi-size Self-attention (MSA) mechanism. MSA splits $\mathcal{C}_{i,j}^*$ into patches of different sizes $p_s \in \{h^j, h^j/2, h^j/4, h^j/8\}$ and computes patch-wise self-

attention [54,55,58,60,62] in different heads. Different heads correspond to different patch sizes, as shown in Fig. 3. To be specific, we first extract patches of shape $2c^j \times p_s \times p_s$ from $\mathcal{C}_{i,j}^*$, and flatten them into 1-dimension vectors for the s -th head. And then we use fully-connected layers to embed the flattened vectors into query embeddings $\mathcal{Q}_{i,j}^s \in \mathbb{R}^{\mathcal{N} \times c^s}$, where $\mathcal{N} = (h^j/p_s) \times (w^j/p_s)$ and $c^s = 2c^j \times p_s \times p_s$. We obtain key embeddings $\mathcal{K}_{i,j}^s$ and value embeddings $\mathcal{V}_{i,j}^s$ with the similar operations. The attention matrix is calculated by the following process:

$$\mathcal{A}_{i,j}^s = \text{softmax} \left(\frac{\mathcal{Q}_{i,j}^s (\mathcal{K}_{i,j}^s)^T}{c^s} \right) \mathcal{V}_{i,j}^s \quad (3)$$

After that, $\mathcal{A}_{i,j}^s$ is reshaped to the original spatial resolution. Similar operations are implemented to obtain features from heads of different patch sizes. Finally, these features are concatenated and passed through a 2D residual block to obtain the output $\mathcal{T}_{i,j} \in \mathbb{R}^{2c^j \times h^j \times w^j}$. We stack the MSA for N times ($N = 3$ in this paper). To further fuse multiple scales of $\mathcal{T}_{i,j}$, we use another MF block to obtain $\mathcal{T}_{i,j}^*$, which is the output of the skip-connection as shown in Fig. 3.

3.4. Anomaly Generation Strategies

To alleviate the data imbalance problem, we propose two kinds of online anomaly generation strategies that can generate various types of anomalies. One strategy is to create

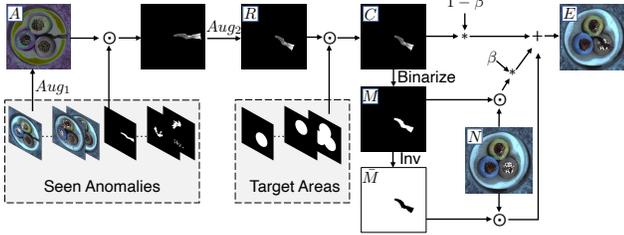


Figure 5. Generating extended anomalies. The anomalous region is augmented by random augmentation and placed on a target area of the normal sample to generate various anomalies online.

in-distribution anomalies by placing augmented anomalous regions from seen anomalies on normal samples, and these generated anomalies are named extended anomalies (EA). EA enlarge the amount of anomalies and mitigate the data imbalance problem. Another strategy is to create out-of-distribution anomalies [68] using normal samples without knowledge of the seen anomalies. These generated anomalies are named simulated anomalies (SA), which supplement potential unseen anomalies.

Extended Anomalies. Instead of simply augmenting the entire image from the seen anomalies, we augment the specific anomalous regions of the seen anomalies and place them at any possible position within the normal sample. First, augmentations (Fig. 5, Aug_1) are applied to a randomly selected anomaly from the seen anomalies in order to generate color varieties (Fig. 5, A). Aug_1 takes two random operations from { equalize, solarize, posterize, sharpness, auto-contrast, invert, gamma-contrast }. After that, we augment the selected anomaly with random spatial transformations as { rotate, shear, shift } to obtain position and shape varieties (Fig. 5, R). Since the extended anomalies should be as realistic and reasonable as possible, we propose a soft position constrain to place R in the foreground. More specifically, Target Areas (TA) is used to refer to areas where anomalies can be placed. We crop R , using a randomly sampled target area, to obtain clipped anomaly region (Fig. 5, C). If R has no overlap with the target region, we perform Aug_2 again until R has overlap with the target region. We binarized C to obtain the ground truth mask (Fig. 5, M). The proposed extended anomalies (Fig. 5, E) is therefore defined as:

$$E = \bar{M} \odot N + (1 - \beta)C + \beta (M \odot N) \quad (4)$$

where \bar{M} is the inverse of M , \odot is the element-wise multiplication operation, β is the opacity parameter [68] for better combination of abnormal and normal parts. For object datasets and texture datasets, the target areas are part of the foreground of the object and part of the whole image, respectively. The shapes of the target are the set of geometries: {circle, rectangular, polygonal}.

Simulated Anomalies. Similar to DRAEM [68], we

multiply Perlin [37] noise with random textures from the DTD [8] dataset and apply these augmented textures to normal images. As these anomalies significantly differs from the seen anomalies, we refer to these out-of-distribution anomalies as heterologous anomalies (HEA). To further expand the diversity of simulated anomalies, we introduce homologous anomalies (HOA), in which anomalies multiplied by the Perlin noise are augmented normal images. Note that the TA mentioned above is also applied to the generation of simulated anomalies. More details about HEA and HOA are presented in the supplementary materials.

3.5. Training and Inference

The decoder of PRN outputs an anomaly score map \mathcal{M}_o , which is of the same shape as the ground truth mask \mathcal{M} . Inspired by [68] and [63], a focal loss [29] and a smooth L1 loss [14] are applied to increase the robustness toward accurate segmentation of hard examples and reduce the over-sensitivity to outliers, respectively. Thus, the total loss \mathcal{L}_{total} used for training PRN is defined as

$$\mathcal{L}_{total} = \text{Smooth}_{\mathcal{L}1}(\mathcal{M}_o, \mathcal{M}) + \lambda \mathcal{L}_{focal}(\mathcal{M}_o, \mathcal{M}) \quad (5)$$

When the predicted \mathcal{M}_o is accurate and sufficiently close to \mathcal{M} , \mathcal{M}_o can be interpreted not only as the pixel-level anomaly localization result, but also as an image-level anomaly estimation for anomaly detection. Specifically, we take the average of the top-K anomalous pixels as the image-level anomaly score for anomaly detection. In a preliminary study, we trained a classification network based on \mathcal{M}_o for image-level anomaly detection, but did not observe an improvement over top-K estimation.

4. Experiments

4.1. Experimental Details

Datasets. To validate the effectiveness and generalizability of our approach, we perform experiments on various datasets, *i.e.*, MVTec Anomaly Detection (MVTec AD [4]), DAGM [61], BeanTech anomaly detection dataset (BTAD [34]), and KolektorSDD2 [7]. There are 10 object sub-datasets and 5 texture sub-datasets in MVTec AD. Each sub-dataset presents a diverse set of anomalies, which enables a general evaluation of surface anomaly detection methods. DAGM contains 10 textured objects with small abnormal regions that are visually very similar to the background. BTAD includes three categories of real-world industrial products showcasing different body and surface defects. KolektorSDD2 is a dataset of surface defects that vary in shape, size, and color, from small scratches and spots to large surface defects. We adopt the general supervised setting [13, 35], where the training set of each sub-dataset contains only 10 abnormal samples. More details will be provided in the supplementary materials.

Category	Unsupervised							Supervised		
	KDAD [47]	CFLOW [18]	DRAEM [68]	SSPCAB [39]	CFA [25]	RD4AD [12]	PatchCore [41]	DevNet [35]	DRA [13]	Ours
Carpet	80.3/95.5	97.6/ 99.2	96.9/97.5	93.1/92.6	99.3/98.6	98.7/98.9	99.1/99.0	82.5/97.2	92.5/98.2	99.7/99.0
Grid	75.3/89.4	98.1/98.9	99.9/ 99.7	99.7/99.5	98.6/97.6	100/98.3	97.3/98.7	90.6/87.9	98.6/86.0	99.4/98.4
Leather	92.3/98.1	99.9/99.7	100/99.0	98.7/96.3	100/99.1	100/99.4	100/99.3	92.2/94.2	98.9/93.8	100/99.7
Tile	91.5/80.2	97.1/96.2	100/99.2	100/99.4	99.2/95.1	99.7/95.7	99.3/95.8	99.9/92.7	100/92.3	100/99.6
Wood	94.5/85.3	98.7/86.0	99.5/95.5	98.4/96.5	100/94.7	99.5/95.8	99.6/95.1	97.9/86.4	99.1/82.9	100/97.8
Bottle	99.2/95.7	99.9/97.2	98.0/99.1	95.6/99.2	100/98.6	100/98.8	100/98.6	99.7/93.9	100/91.3	100/99.4
Cable	90.3/80.2	97.6/97.8	90.9/95.2	92.7/95.1	99.9/98.8	96.1/97.2	99.9/98.5	98.7/88.8	94.2/86.6	98.9/ 98.8
Capsule	81.4/95.2	97.0/ 99.1	91.3/88.1	96.9/90.2	97.4/98.4	96.1/98.7	98.0/99.0	71.9/91.8	95.1/89.3	98.0/98.5
Hazelnut	98.8/95.0	100/98.8	100/99.7	100/99.7	100/98.6	100/99.0	100/98.7	99.7/91.1	100/89.6	100/99.7
Metal Nut	77.1/83.3	98.5/98.6	100/99.6	100/99.4	100/98.7	100/97.3	99.9/98.3	98.8/77.8	99.1/79.5	100/99.7
Pill	84.4/89.9	96.2/98.9	97.1/97.3	97.4/97.2	97.7/98.0	98.7/98.1	97.5/97.6	87.1/82.6	88.3/84.5	99.3/99.5
Screw	82.4/95.8	93.1/98.9	98.7/99.3	97.8/99.0	95.1/98.3	97.8/ 99.7	98.2/99.5	97.2/60.3	99.5/54.0	95.9/97.5
Toothbrush	97.1/95.5	98.8/99.0	100/97.3	97.9/97.3	100/98.8	100/99.1	100/98.6	79.2/84.6	87.5/75.5	100/99.6
Transistor	84.9/75.9	92.9/98.2	91.7/85.2	88.0/84.8	100/98.1	95.5/92.3	99.9/96.5	89.1/56.0	88.3/79.1	99.7/ 98.4
Zipper	93.7/95.3	97.1/ 99.1	100/99.1	100/98.4	99.5/98.6	97.9/98.3	99.5/98.9	99.1/93.7	99.7/96.9	99.7/98.8
Average	88.2/90.0	97.5/97.7	97.6/96.7	97.1/96.3	99.1/98.0	98.7/97.8	99.2/98.1	92.2/85.3	96.1/85.3	99.4/99.0

Table 1. Anomaly Detection and Localization on MVTec [4]. Best results on Image AUROC or Pixel AUROC are highlighted in bold.

Evaluation Metrics. Following previous work, we evaluate the results via the area under the receiver operating characteristic curve at the image level (Image-AUROC) and pixel level (Pixel-AUROC). However, anomalous regions typically only occupy a tiny fraction of the entire image. Thus, Pixel-AUROC can not accurately reflect the localization accuracy due to the fact that the false positive rate is dominated by the extremely high number of non-anomalous pixels and remains low despite false positive detection [53]. To comprehensively measure localization performance, we introduce Per Region Overlap (PRO) [5] score and pixel-level Average Precision (AP) [68]. The PRO score treats anomaly regions of varying sizes equally [12, 41], while AP is more appropriate for highly imbalanced classes, especially for industrial anomaly localization, where accuracy is critical [68].

Implementation Details. All images in four datasets are resized to 256×256 . We use layer1, layer2 and layer3 of ResNet-18 [21] pre-trained on ImageNet to obtain feature maps with $64 \times 64 \times 64$, $128 \times 32 \times 32$ and $256 \times 16 \times 16$ scales respectively and frozen these blocks during training. The number of prototypes depends on the dataset and accounts for 10% of the total number of normal samples in the dataset. The maximum number of iterations of k-means for each scale is set to 300. We use Adam optimizer [24] for the parameter optimization, with an initial learning rate 10^{-4} and a weight decay of 10^{-2} . α and γ in focal loss is set to 0.5 and 4 respectively. λ in the total loss is set to 5. We train for 700 epochs with a batch size of 64 consisting of 32 normal samples, 16 extended anomalies and 16 simulated anomalies to ensure the diversity of anomalies. We take the average of the top 100 anomalous pixels as the image-level anomaly score. We compare PRN to seven unsupervised SOTA methods and two supervised SOTA methods. The results we report are based on the implementation provided by these methods. The backbone of PatchCore [41],

RD4AD [12], CFLOW [18] and CFA [25] is WideResNet50. SSPCAB [1, 39] replaces the penultimate convolutional layer of reconstructive encoder in DRAEM [68]. DevNet [35] proposes that the anomaly score given by the network can be further back-propagated to the original image pixels to infer which pixels are the major contributors to the anomaly for anomaly localization. We use this approach to obtain the anomaly localization performance of DRA [13].

4.2. Anomaly Detection and Localization on MVTec

Anomaly detection and localization results on MVTec are shown in Table 1. Our method achieves the highest image AUROC (detection) and the highest pixel AUROC (localization) in 10 out of 15 classes. The average image AUROC results show that our method outperforms unsupervised SOTA by 0.2% and supervised SOTA by 3.3%. Meanwhile, for pixel AUROC, our method outperforms unsupervised SOTA by 0.9% and supervised SOTA by 13.7%.

For a comprehensive presentation of the capabilities on anomaly localization, two additional metric results, PRO and AP, are shown in Table 2. PRN outperforms the previous unsupervised SOTA by 2.2% and the previous supervised SOTA by 22.8% on the PRO metric. This confirms that PRN is more effective at simultaneously localizing anomalous regions of varying sizes. The more challenging AP metric further demonstrates the excellent anomaly localization capability of PRN. A better AP score is achieved in 12 out of 15 classes and is comparable to SOTA in other classes. In terms of overall AP, our approach even outperforms unsupervised SOTA by 10.5% and supervised SOTA by 52.6%. This significant improvement over AP goes a long way to demonstrate that PRN is more discriminative between normal and abnormal pixels. We further compare the pre-trained based approaches in terms of inference time per image (second) and performance, as shown in Table 3. All experiments were conducted on NVIDIA GeForce RTX 3090, using a uniform

Category	Unsupervised							Supervised		
	KDAD [47]	CFLOW [18]	DRAEM [68]	SSPCAB [39]	CFA [25]	RD4AD [12]	PatchCore [41]	DevNet [35]	DRA [13]	Ours
Carpet	92.5/45.6	97.6 /68.3	92.9/65.1	86.4/48.6	93.6/57.2	95.4/56.5	95.5/62.2	85.8/45.7	92.2/52.3	97.0 / 82.0
Grid	72.9/7.3	96.0/41.2	98.3 / 62.8	98.0/57.9	92.9/25.8	94.2/15.8	94.0/24.5	79.8/25.5	71.5/26.8	95.9/45.7
Leather	97.5/26.8	99.2/64.5	97.4/ 72.9	94.0/60.7	95.4/48.5	98.2/47.6	96.9/45.3	88.5/8.1	84.0/5.6	99.2 /69.7
Tile	74.3/27.7	89.1/60.1	98.2 /95.2	98.1/96.1	83.3/55.9	85.6/54.1	91.3/56.2	78.9/52.3	81.5/57.6	98.2 / 96.5
Wood	76.5/24.3	82.8/29.0	90.3/74.6	92.8/78.9	85.9/49.0	91.4/48.3	87.1/49.3	75.4/25.1	69.7/22.7	95.9 / 82.6
Bottle	88.6/54.8	94.0/68.1	96.8/88.9	96.3/89.4	94.6/80.3	96.3/78.0	95.4/76.8	83.5/51.5	77.6/41.2	97.0 / 92.3
Cable	66.2/12.6	94.1/60.6	81.0/56.4	80.4/52.0	91.7/74.7	94.1/52.6	96.8/67.0	80.9/36.0	77.7/34.7	97.2 / 78.9
Capsule	90.1/10.1	94.0/48.8	82.7/39.6	92.5/46.4	93.0/48.3	95.5 /47.2	93.4/46.0	83.6/15.5	79.1/11.7	92.5/ 62.2
Hazelnut	94.3/34.2	97.1/59.9	98.5 /92.6	98.2/93.4	95.2/60.0	96.9/60.7	90.9/53.2	83.6/22.1	86.9/22.5	97.4/ 93.8
Metal Nut	76.9/34.1	91.5/88.0	97.0/97.0	97.7 /94.7	91.4/92.2	94.9/78.6	92.6/86.6	76.9/35.6	76.7/29.9	95.8/ 98.0
Pill	86.4/20.9	95.2/82.0	88.4/47.6	89.6/48.3	95.4/81.9	96.7/76.5	94.5/75.7	69.2/14.6	77.0/21.6	97.2 / 91.3
Screw	85.2/6.1	95.8/43.9	95.0/ 66.5	95.2/61.7	93.5/28.7	98.5 /52.1	97.5/34.7	31.1/1.4	30.1/5.0	92.4/44.9
Toothbrush	87.3/18.3	95.3/46.3	85.6/45.5	85.5/39.3	86.8/55.7	92.3/51.1	94.0/37.9	33.5/6.7	56.1/4.5	95.6 / 78.1
Transistor	68.1/25.8	82.5/67.5	70.4/39.0	62.5/38.1	95.1 /76.2	83.3/54.1	92.3/66.9	39.1/6.4	49.0/11.0	94.8/ 85.6
Zipper	86.5/31.5	96.6/65.2	96.8 / 77.6	95.2/76.4	94.3/65.2	95.3/57.5	96.1/62.3	81.3/19.6	91.0/42.9	95.5/ 77.6
Total Average	82.9/25.34	93.4/59.6	91.3/68.1	90.8/65.5	92.1/60.0	93.9/55.4	93.9/56.3	71.4/24.4	73.3/26.0	96.1 / 78.6

Table 2. Results of the PRO and AP metrics for anomaly localization performance on MVTEC [4].

	Backbone	I \uparrow	P \uparrow	O \uparrow	A \uparrow	T \downarrow
CFLOW	WResNet50	97.5	97.7	93.4	59.6	0.127
RD4AD		98.7	97.8	93.9	55.4	0.094
PatchCore		99.2	98.1	93.9	56.3	0.133
CFLOW	ResNet18	96.2	98.1	92.8	59.2	0.106
RD4AD		97.9	97.1	92.7	53.7	0.076
DRA		96.1	84.1	71.5	25.7	0.223
PRN(Ours)		99.4	99.0	96.1	78.6	0.064

Table 3. Comparison of pre-trained based approaches in terms of performance and inference time (second) on MVTEC [4]. “I”, “P”, “O”, “A” and “T” respectively refer to the five metrics of image auROC, pixel auROC, pixel pro, pixel ap, and inference time per image.

standard. Our approach not only gains the best performance, but also significantly reduces the inference time.

We qualitatively evaluate the performance of anomaly localization compared to state-of-the-art methods DRAEM [68] and PatchCore [41] by visualizing the results in Fig. 6. Our model accurately locates the anomalies and clearly focus on all anomalous regions, regardless of their sizes, shapes and numbers. Additional qualitative results are provided in the supplementary material.

4.3. Ablation Study

The importance of MP, MSA and MF. We investigate the importance of each modules in PRN and the results are reported in Table 4. We have the U-Net-like architecture without any module on the skip-connection branch as the baseline. Overall, PRN outperforms the baseline by a large margin, especially on the P, O, and A metrics. All metrics are significantly boosted by employing the MP that performs explicit residual representation. When applying the MSA which performs variable-sized anomalous feature learning, the performance is further improved. This confirms the effectiveness of information exchanging across multi-size receptive fields. Finally, removing the MF causing the degradation

U-Net	Module				Performance			
	MP	MSA	MF	I \uparrow	P \uparrow	O \uparrow	A \uparrow	
\checkmark				97.4	91.7	88.6	58.5	
\checkmark	\checkmark			98.9	98.5	95.3	77.0	
\checkmark		\checkmark		97.8	97.0	92.1	74.0	
\checkmark	\checkmark	\checkmark		98.7	98.5	95.4	78.1	
\checkmark	\checkmark	\checkmark	\checkmark	99.4	99.0	96.1	78.6	

Table 4. Ablations of different modules in PRN.

of performance, indicates that it is necessary to exchange information across different scales.

Effects of different anomaly generation strategies. We perform ablation studies to investigate the impact of the different components of the proposed anomaly generation strategies in Table 5. The proposed EA alleviates the problem of seen appearance variance, but does not adequately explore the underlying unseen anomalies. Table 5 indicates that the performance of the model increases with the variety of generated anomalies. We argue that the proposed SA consisting of both HEA and HOA can generate anomalies of various sizes, shapes and numbers, allowing our model to generalize to unseen anomalies. Besides, the proposed TA imposes soft constraints on the locations where anomaly regions are imposed, making the generated anomalies as realistic and reasonable as possible, thus significantly improving the performance of the model.

The effect of prototype proportion. The effect of the ratio of prototypes to total normal samples is compared in Table 6. Note that 100% means that no clustering is performed. Each feature map of a normal sample is regarded as a prototype and the number of prototypes is equal to the number of normal samples. The poor performance of the PRN_{100%} indicates that the residual representation obtained from the closest cluster prototype is more representative than that obtained from the single closest sample. Besides, too

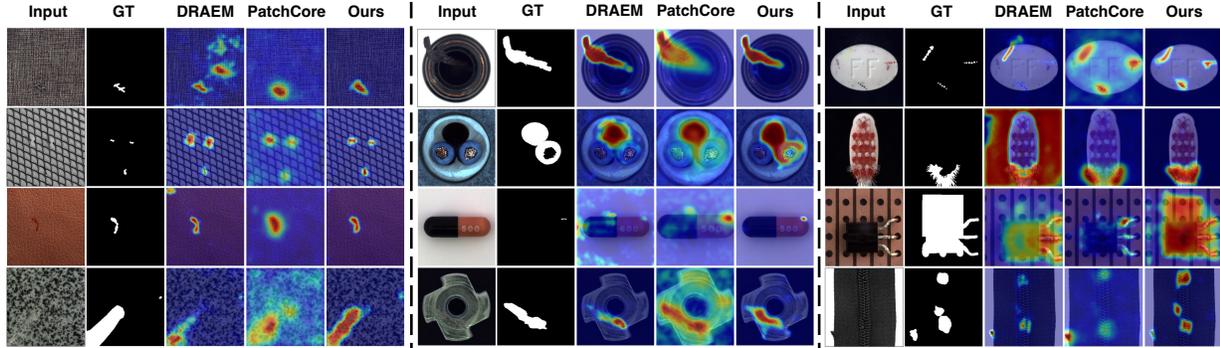


Figure 6. Qualitative examples on MVTeC [4]. PRN achieves more accurate localization results for various types of anomalies.

Anomaly Generation				Performance			
EA	HEA	HOA	TA	I \uparrow	P \uparrow	O \uparrow	A \uparrow
✓			✓	98.6	97.2	93.4	75.7
✓	✓		✓	99.1	98.4	95.4	77.4
✓		✓	✓	98.6	98.4	95.7	75.2
	✓	✓	✓	98.7	98.2	95.1	73.4
✓	✓	✓		98.4	98.4	94.9	77.6
✓	✓	✓	✓	99.4	99.0	96.1	78.6

Table 5. Ablations of anomaly generation strategies.

	I \uparrow	P \uparrow	O \uparrow	A \uparrow	T \downarrow
PRN _{5%}	99.2	98.6	95.4	78.1	0.063
PRN _{10%}	99.4	99.0	96.1	78.6	0.064
PRN _{20%}	99.2	98.8	95.7	77.3	0.066
PRN _{100%}	86.2	91.4	75.4	49.9	0.074

Table 6. Ablations of the ratio of prototypes to total normal samples.

	DevNet [35]				DRA [13]				PRN(Ours)			
	I \uparrow	P \uparrow	O \uparrow	A \uparrow	I \uparrow	P \uparrow	O \uparrow	A \uparrow	I \uparrow	P \uparrow	O \uparrow	A \uparrow
1	79.6	75.3	51.0	16.5	88.9	78.8	58.2	19.1	98.8	98.3	95.4	74.7
5	86.7	83.7	66.9	22.7	93.5	82.8	68.6	21.9	99.2	98.6	95.6	76.4
10	92.2	85.3	71.4	24.4	96.1	85.3	73.3	26.0	99.4	99.0	96.1	78.6

Table 7. Impact of the number of seen anomalies used.

few prototypes lead to insufficient discrimination between prototypes, resulting in inferior performance. The results indicate that a proportion of 10% produces the optimum performance. In addition, using fewer prototypes can speed up inference.

Effects of the number of seen anomalies used. As shown in Table 7, we explore the impact of the number of anomalies used. Our approach significantly outperforms Devnet [35] and DRA [13] using different numbers of seen anomalies, which demonstrates the effectiveness of our proposed anomaly generation strategies and the robustness of PRN to datasets of different levels of imbalance.

4.4. Evaluation on other benchmarks

To further evaluate the anomaly detection and localization capabilities of PRN, we benchmark PRN on three additional

	DAGM [61]				BTAD [34]				KolektorSDD2 [7]			
	I \uparrow	P \uparrow	O \uparrow	A \uparrow	I \uparrow	P \uparrow	O \uparrow	A \uparrow	I \uparrow	P \uparrow	O \uparrow	A \uparrow
DRAEM	91.1	83.4	70.5	35.6	89.0	87.1	61.6	19.2	81.1	85.6	67.9	39.1
CFLOW	91.2	95.1	87.6	45.2	90.5	96.1	71.6	54.0	95.2	97.4	93.8	46.0
SSPCAB	90.4	84.5	71.9	33.9	88.3	83.5	54.1	13.0	83.4	86.2	66.1	44.5
RD4AD	90.7	94.1	85.5	40.8	94.4	96.9	75.8	53.5	96.0	97.6	94.7	43.5
PatchCore	92.5	96.1	88.0	49.0	92.6	96.9	76.3	51.5	94.6	97.1	89.3	49.8
DRA	93.5	95.1	88.8	47.6	94.2	75.4	56.2	12.4	86.8	84.4	56.9	3.6
Ours	98.2	96.6	93.8	49.4	94.7	97.1	78.0	54.0	96.4	97.6	94.9	72.5

Table 8. Comparison of PRN with other approaches on DAGM, BTAD, and KolektorSDD2.

widely used datasets, namely DAGM [61], BTAD [34] and KolektorSDD2 [7]. As shown in Table 8, PRN achieves new SOTA performance on all three datasets, proving its effectiveness and generalization. Results for more detailed comparisons and some qualitative examples are provided in the supplementary material.

5. Conclusion

In this paper, we proposed a novel framework called Prototypical Residual Network for anomaly detection and localization. PRN learns residual representations across multi-scale feature maps and within multi-size receptive fields at each scale, enabling accurate detection and localization of anomalous regions that come in a variety of sizes, shapes and numbers. In addition, we propose various anomaly generation strategies to expand and diversify the anomalies. We conduct in-depth experiments on four popular datasets to confirm the effectiveness and generalizability of our approach. PRN achieves new SOTA on anomaly detection and significantly surpasses previous arts in anomaly localization performance.

Limitations. Our approach requires the dataset to provide accurate ground truth masks for anomalies. Using a single image-level anomaly average score for anomalous images with different defect sizes does not favor tiny defects. We leave this intriguing extension to future work.

Acknowledgement This project was supported by NSFC under Grant No. 62102092 and No. 62032006.

A. Appendix

A.1. Anomaly Generation Strategies

This section details the generation of simulated anomalies, as shown in Fig. 7. A noise image is generated by a Perlin noise generator [37, 68] (Fig. 7, P), and then the noise parts within a target area are retained as the ground truth mask (Fig. 7, M). As the shape, size, and number of generated anomalous regions vary widely, we synthesize simulated anomalies (Fig. 7, S) as:

$$S = \bar{M} \odot N + (1 - \beta)(M \odot A) + \beta(M \odot N) \quad (6)$$

where N is the normal sample, A is the source image of the anomaly, \bar{M} is the inverse of M , \odot is the element-wise multiplication operation, β is the opacity parameter for better combination of abnormal and normal regions. When A is an image randomly sampled from the DTD dataset [8] and is augmented (A_{ug_1} , Fig. 5 in Section 3.4), we define S as a HETerologous Anomaly (HEA). Correspondingly, when A is an image randomly sampled from augmented normal samples, we define S as a HOMology Anomaly (HOA). In particular, the normal image is first augmented (A_{ug_1} , Fig. 5 in Section 3.4), then is evenly divided into an 8×8 grid and randomly arranged before being reassembled [63].

Fig. 8 shows the anomalies generated by different strategies. In addition to increasing the number, extended anomalies (EA) increase the variety of seen anomalies. HEA and HOA supplement potential unseen anomalies with anomalies significantly different from seen anomalies.

A.2. Dataset Split

MVTec AD [4] is a widely used anomaly detection and localization benchmark with 15 classes, each containing one to several subclasses of anomalies. Following the general setting proposed by DRA [13], the 10 labeled anomaly samples are sampled from all possible anomaly classes in the test set per dataset. These sampled anomalies are then removed from the test data. Both BTAD [34] and KolektorSDD2 [7] are real-world industrial datasets containing three product types and one product type, respectively. The general setting used in BTAD and SDD2 is same to that used in MVTec. DAGM [61] contains 10 texture classes, and the original training set for each class consists of normal and abnormal samples. For each class, we first move all anomalous samples from the original training set to the original test set, and then randomly select ten anomalous samples from the test set as part of the new training set. These sampled anomalies are then removed from the test set.

A.3. More Detailed Comparison

Table 9 includes fine-grained anomaly detection and localization performance comparisons on all DAGM sub-datasets. We observe that PRN consistently performs well on all 10

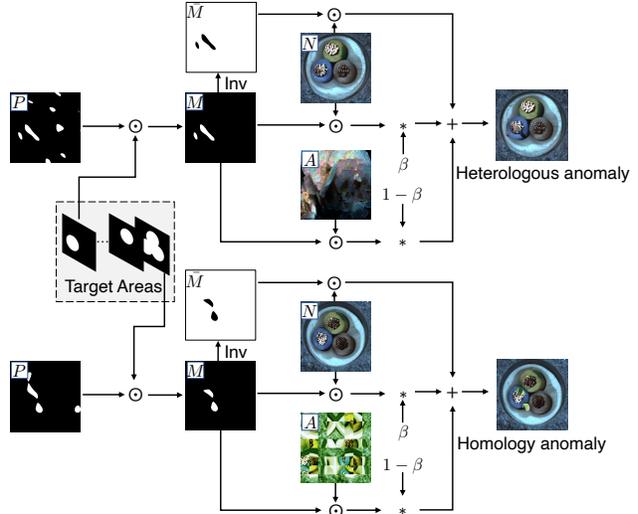


Figure 7. Generating simulated anomalies.

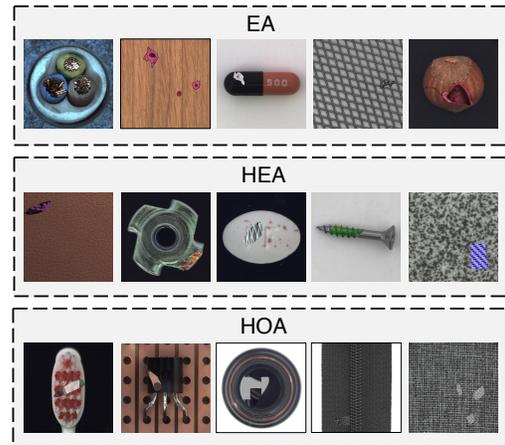


Figure 8. Examples of anomalies generated by different strategies.

sub-datasets and, in the average scenario, performs best across all four criteria. In particular, our approach outperforms previous methods by a large margin in two metrics, image auROC and pro.

We also compare the anomaly detection and location performance of each method in detail on the three BTAD products and report the numerical results in Table 10. It can be concluded that our method achieves consistently higher performance than the others on different categories .

A.4. More Qualitative Examples

We further qualitatively evaluate the performance of anomaly detection and location compared to state-of-the-art methods by introducing additional visualizations, as shown in Fig. 9, Fig. 10 and Fig. 11. Our method accurately detects and localizes anomalies in a wide range of sizes, shapes and

Category	DRAEM [68]				CFLOW [18]				SSPCAB [39]				RD4AD [12]				PatchCore [41]				Ours			
	I↑	P↑	O↑	A↑	I↑	P↑	O↑	A↑	I↑	P↑	O↑	A↑	I↑	P↑	O↑	A↑	I↑	P↑	O↑	A↑	I↑	P↑	O↑	A↑
Class1	86.9	75.4	56.3	20.9	91.6	94.1	84.1	33.4	95.3	78.9	61.2	29.9	95.2	92.8	83.0	41.6	84.4	89.6	72.8	13.1	100	92.7	90.3	50.1
Class2	85.8	83.7	66.1	18.2	98.2	99.6	98.2	50.2	93.9	92.0	80.4	18.3	99.7	99.7	99.1	57.8	100	99.7	99.3	55.5	96.0	97.1	95.6	44.8
Class3	98.0	90.3	78.2	32.9	88.3	93.7	86.0	32.9	99.6	90.3	79.2	31.7	81.2	93.9	85.2	31.8	94.0	96.2	92.4	50.9	99.2	94.2	91.3	32.4
Class4	99.3	98.6	95.5	62.4	100	99.5	98.5	65.1	99.9	99.1	97.6	74.7	99.9	99.1	97.7	64.6	100	99.4	98.4	88.2	99.7	98.2	96.7	67.2
Class5	97.9	56.4	39.9	21.9	86.3	94.3	84.5	50.7	81.1	53.6	35.9	15.5	74.1	86.7	64.3	31.2	90.6	95.2	77.3	29.6	96.9	94.9	86.1	30.2
Class6	100	96.0	89.3	71.5	96.5	96.1	87.9	46.9	100	95.4	88.3	70.0	92.0	88.3	68.9	30.3	99.4	98.1	93.5	71.2	100	98.4	95.7	71.7
Class7	100	96.7	90.8	58.1	98.9	96.0	91.8	61.4	100	94.8	87.0	51.1	99.8	95.2	91.4	65.7	99.9	96.9	94.8	77.7	100	95.1	91.3	51.3
Class8	99.7	92.9	90.4	34.2	56.7	79.9	51.0	3.2	96.4	91.1	88.9	23.2	65.2	86.2	67.6	7.0	60.6	86.4	56.5	7.8	93.4	97.1	95.1	34.4
Class9	50.2	49.7	13.3	0.1	99.9	99.9	99.8	65.1	50.9	60.4	26.1	0.1	100	99.8	99.4	26.5	96.4	99.4	95.7	45.9	97.1	98.7	96.8	46.4
Class10	92.7	94.2	85.4	35.7	95.7	98.0	94.4	42.9	86.5	89.1	74.7	24.4	99.6	99.0	97.9	51.1	99.9	99.6	99.0	49.6	99.9	99.6	99.0	65.6
Average	91.1	83.4	70.5	35.6	91.2	95.1	87.6	45.2	90.4	84.5	71.9	33.9	90.7	94.1	85.5	40.8	92.5	96.1	88.0	49.0	98.2	96.6	93.8	49.4

Table 9. Anomaly Detection and Localization on DAGM [61]. “I”, “P”, “O” and “A” respectively refer to the five metrics of image auROC, pixel auROC, pro and ap. The best results are highlighted in bold.

Category	DRAEM [68]				CFLOW [18]				SSPCAB [39]				RD4AD [12]				PatchCore [41]				Ours			
	I↑	P↑	O↑	A↑	I↑	P↑	O↑	A↑	I↑	P↑	O↑	A↑	I↑	P↑	O↑	A↑	I↑	P↑	O↑	A↑	I↑	P↑	O↑	A↑
01	98.5	91.5	61.4	17.0	93.4	94.8	60.1	39.6	96.2	92.4	62.8	18.1	98.8	95.7	72.8	49.3	96.6	96.5	78.4	47.1	100	96.6	81.4	38.8
02	68.6	73.4	39.0	23.3	79.0	93.9	56.9	65.5	69.3	65.6	28.6	15.8	84.9	96.0	55.8	66.1	81.3	94.9	54.0	56.3	84.1	95.1	54.4	65.7
03	99.8	96.3	84.3	17.2	99.1	99.5	97.9	56.8	99.4	92.4	71.0	5.0	99.5	99.0	98.8	45.1	99.9	99.2	96.4	51.2	99.9	99.6	98.3	57.4
Average	89.0	87.1	61.6	19.2	90.5	96.1	71.6	54.0	88.3	83.5	54.1	13	94.4	96.9	75.8	53.5	92.6	96.9	76.3	51.5	94.7	97.1	78.0	54.0

Table 10. Anomaly Detection and Localization on BTAD [34].

numbers, as demonstrated by qualitative comparison results. Moreover, we argue that some of the localization errors can be attributed to inaccurate ground truth labels on anomalies. An example of this is shown in the second row of Fig. 11, where the ground truth does not label all anomalous regions. Another example is shown on the left in the fourth row of Fig. 10, where the ground truth labels a broad anomaly region, but our method correctly localizes the anomaly region. These imprecise annotations inevitably impact the anomaly localization scores of the evaluated methods.

References

[1] Samet Akcay, Dick Ameln, Ashwin Vaidya, Barath Lakshmanan, Nilesh Ahuja, and Utku Genc. Anomalib: A

deep learning library for anomaly detection. *arXiv preprint arXiv:2202.08341*, 2022. 6

- [2] Samet Akcay, Amir Atapour-Abarghouei, and Toby P Breckon. Ganomaly: Semi-supervised anomaly detection via adversarial training. In *ACCV*, 2018. 3
- [3] Samet Akçay, Amir Atapour-Abarghouei, and Toby P Breckon. Skip-ganomaly: Skip connected and adversarially trained encoder-decoder anomaly detection. In *IJCNN*, 2019. 2
- [4] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Mvtec ad—a comprehensive real-world dataset for unsupervised anomaly detection. In *CVPR*, 2019. 1, 5, 6, 7, 8, 9, 10
- [5] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Uninformed students: Student-teacher

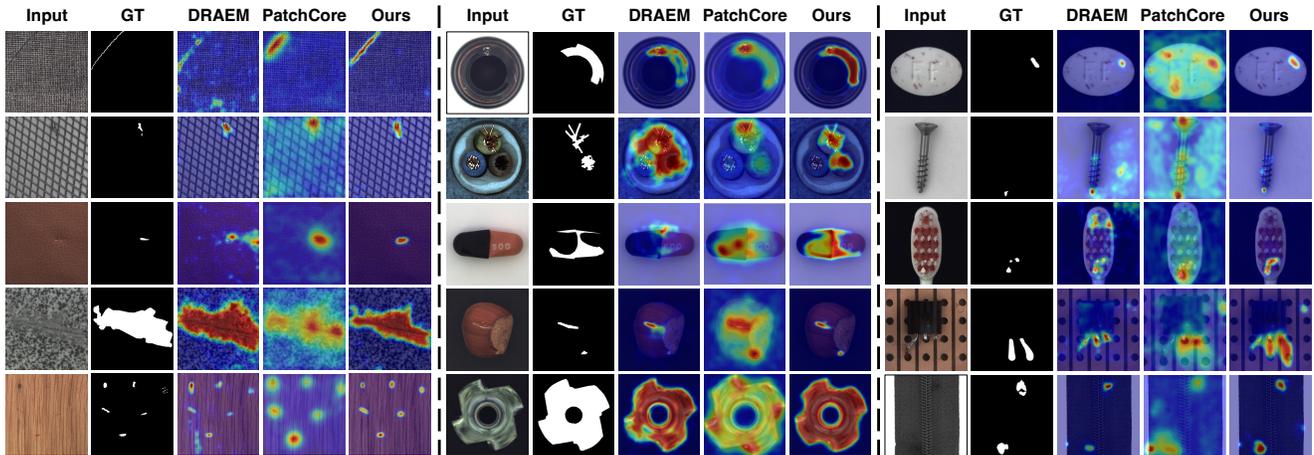


Figure 9. More qualitative examples on MVTec [4].

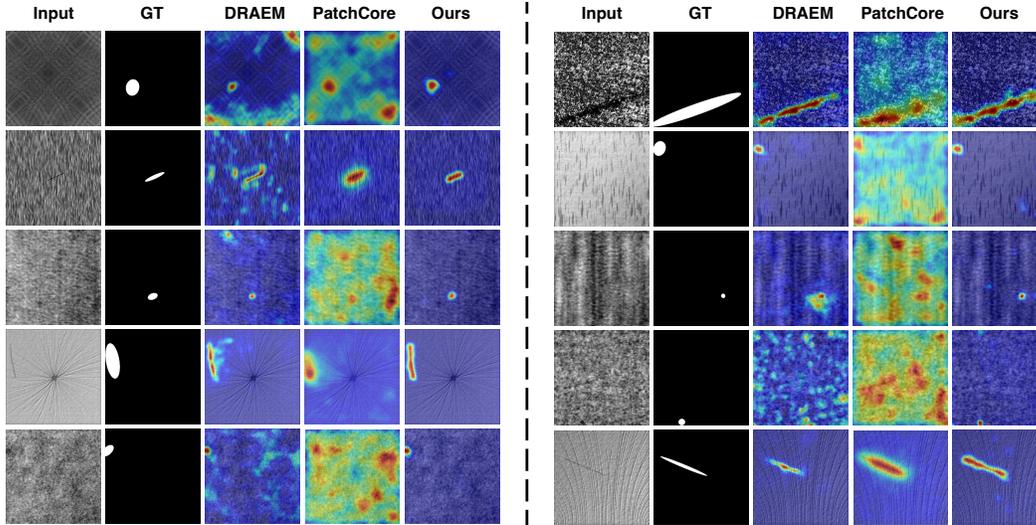


Figure 10. Qualitative examples on DAGM [61].

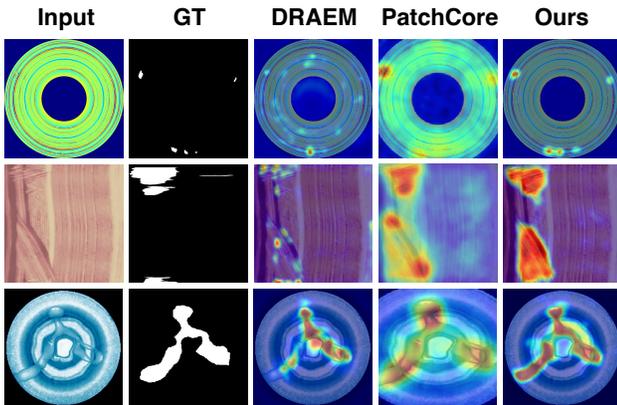


Figure 11. Qualitative examples on BTAD [34].

anomaly detection with discriminative latent embeddings. In *CVPR*, 2020. 2, 3, 6

- [6] Paul Bergmann, Sindy Löwe, Michael Fauser, David Sattlegger, and Carsten Steger. Improving unsupervised defect segmentation by applying structural similarity to autoencoders. *arXiv preprint arXiv:1807.02011*, 2018. 2
- [7] Jakob Božič, Domen Tabernik, and Danijel Škočaj. Mixed supervision for surface-defect detection: From weakly to fully supervised learning. *Comput Ind*, 2021. 1, 5, 8, 9
- [8] Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *CVPR*, 2014. 5, 9
- [9] Niv Cohen and Yedid Hoshen. Sub-image anomaly detection with deep pyramid correspondences. *arXiv preprint arXiv:2005.02357*, 2020. 2, 3
- [10] Thomas Defard, Aleksandr Setkov, Angélique Loesch, and Romaric Audigier. Padim: a patch distribution modeling

framework for anomaly detection and localization. In *ICPR*, 2021. 2, 3

- [11] David Dehaene, Oriol Frigo, Sébastien Combexelle, and Pierre Eline. Iterative energy-based projection on a normal data manifold for anomaly localization. *arXiv preprint arXiv:2002.03734*, 2020. 2
- [12] Hanqiu Deng and Xingyu Li. Anomaly detection via reverse distillation from one-class embedding. In *CVPR*, 2022. 2, 3, 6, 7, 10
- [13] Choubo Ding, Guansong Pang, and Chunhua Shen. Catching both gray and black swans: Open-set supervised anomaly detection. In *CVPR*, 2022. 1, 2, 3, 5, 6, 7, 8, 9
- [14] Ross Girshick. Fast r-cnn. In *ICCV*, 2015. 5
- [15] Dong Gong, Lingqiao Liu, Vuong Le, Budhaditya Saha, Moussa Reda Mansour, Svetha Venkatesh, and Anton van den Hengel. Memorizing normality to detect anomaly: Memory-augmented deep autoencoder for unsupervised anomaly detection. In *ICCV*, 2019. 2
- [16] Nico Görnitz, Marius Kloft, Konrad Rieck, and Ulf Brefeld. Toward supervised anomaly detection. *JAIR*, 2013. 3
- [17] Jiaqi Gu, Hyoukjun Kwon, Dilin Wang, Wei Ye, Meng Li, Yu-Hsin Chen, Liangzhen Lai, Vikas Chandra, and David Z Pan. Multi-scale high-resolution vision transformer for semantic segmentation. In *CVPR*, 2022. 3
- [18] Denis Gudovskiy, Shun Ishizaka, and Kazuki Kozuka. Cflow-ad: Real-time unsupervised anomaly detection with localization via conditional normalizing flows. In *WACV*, 2022. 2, 3, 6, 7, 10
- [19] Songqiao Han, Xiyang Hu, Hailiang Huang, Mingqi Jiang, and Yue Zhao. Adbench: Anomaly detection benchmark. *arXiv preprint arXiv:2206.09426*, 2022. 2
- [20] John A Hartigan and Manchek A Wong. Algorithm as 136: A k-means clustering algorithm. *JSTOR*, 1979. 3
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3, 6

- [22] Jinlei Hou, Yingying Zhang, Qiaoyong Zhong, Di Xie, Shiliang Pu, and Hong Zhou. Divide-and-assemble: Learning block-wise memory for unsupervised anomaly detection. In *ICCV*, 2021. 3
- [23] Jin-Hwa Kim, Do-Hyeong Kim, Saehoon Yi, and Taehoon Lee. Semi-orthogonal embedding for efficient unsupervised anomaly segmentation. *arXiv preprint arXiv:2105.14737*, 2021. 3
- [24] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6
- [25] Sungwook Lee, Seunghyun Lee, and Byung Cheol Song. Cfa: Coupled-hypersphere-based feature adaptation for target-oriented anomaly localization. *ACCESS*, 2022. 2, 3, 6, 7
- [26] Chun-Liang Li, Kihyuk Sohn, Jinsung Yoon, and Tomas Pfister. Cutpaste: Self-supervised learning for anomaly detection and localization. In *CVPR*, 2021. 2
- [27] Jie Li, Xing Xu, Lianli Gao, Zheng Wang, and Jie Shao. Cognitive visual anomaly detection with constrained latent representations for industrial inspection robot. *Appl. Soft Comput.*, 2020. 3
- [28] Yufei Liang, Jiangning Zhang, Shiwei Zhao, Runze Wu, Yong Liu, and Shuwen Pan. Omni-frequency channel-selection representations for unsupervised anomaly detection. *arXiv preprint arXiv:2203.00259*, 2022. 2
- [29] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 5
- [30] Wenqian Liu, Runze Li, Meng Zheng, Srikrishna Karanam, Ziyang Wu, Bir Bhanu, Richard J Radke, and Octavia Camps. Towards visually explaining variational autoencoders. In *CVPR*, 2020. 2
- [31] Wen Liu, Weixin Luo, Zhengxin Li, Peilin Zhao, Shenghua Gao, et al. Margin learning embedded prediction for video anomaly detection with a few anomalies. In *IJCAI*, 2019. 3
- [32] Wen Liu, Weixin Luo, Dongze Lian, and Shenghua Gao. Future frame prediction for anomaly detection—a new baseline. In *CVPR*, 2018. 1
- [33] Lingchen Meng, Hengduo Li, Bor-Chun Chen, Shiyi Lan, Zuxuan Wu, Yu-Gang Jiang, and Ser-Nam Lim. Advavit: Adaptive vision transformers for efficient image recognition. In *CVPR*, 2022. 2
- [34] Pankaj Mishra, Riccardo Verk, Daniele Fornasier, Claudio Piciarelli, and Gian Luca Foresti. Vt-adl: A vision transformer network for image anomaly detection and localization. In *ISIE*, 2021. 1, 3, 5, 8, 9, 10, 11
- [35] Guansong Pang, Choubo Ding, Chunhua Shen, and Anton van den Hengel. Explainable deep few-shot anomaly detection with deviation networks. *arXiv preprint arXiv:2108.00462*, 2021. 2, 3, 5, 6, 7, 8
- [36] Guansong Pang, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. Deep learning for anomaly detection: A review. *CSUR*, 2021. 2
- [37] Ken Perlin. An image synthesizer. *ACM SIGGRAPH*, 1985. 5, 9
- [38] Jonathan Pirnay and Keng Chai. Inpainting transformer for anomaly detection. In *ICIAP*, 2022. 3
- [39] Nicolae-Cătălin Ristea, Neelu Madan, Radu Tudor Ionescu, Kamal Nasrollahi, Fahad Shahbaz Khan, Thomas B Moeslund, and Mubarak Shah. Self-supervised predictive convolutional attentive block for anomaly detection. In *CVPR*, 2022. 3, 6, 7, 10
- [40] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 3
- [41] Karsten Roth, Latha Pemula, Joaquin Zepeda, Bernhard Schölkopf, Thomas Brox, and Peter Gehler. Towards total recall in industrial anomaly detection. In *CVPR*, 2022. 1, 2, 3, 6, 7, 10
- [42] Marco Rudolph, Bastian Wandt, and Bodo Rosenhahn. Same same but different: Semi-supervised defect detection with normalizing flows. In *WACV*, 2021. 2, 3
- [43] Marco Rudolph, Tom Wehrbein, Bodo Rosenhahn, and Bastian Wandt. Fully convolutional cross-scale-flows for image-based defect detection. In *WACV*, 2022. 2, 3
- [44] Lukas Ruff, Jacob R Kauffmann, Robert A Vandermeulen, Grégoire Montavon, Wojciech Samek, Marius Kloft, Thomas G Dietterich, and Klaus-Robert Müller. A unifying review of deep and shallow anomaly detection. *IEEE*, 2021. 2
- [45] Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoab Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. Deep one-class classification. In *PMLR*, 2018. 3
- [46] Lukas Ruff, Robert A Vandermeulen, Nico Görnitz, Alexander Binder, Emmanuel Müller, Klaus-Robert Müller, and Marius Kloft. Deep semi-supervised anomaly detection. *arXiv preprint arXiv:1906.02694*, 2019. 2, 3
- [47] Mohammadreza Salehi, Niousha Sadjadi, Soroosh Baselizadeh, Mohammad H Rohban, and Hamid R Rabiee. Multiresolution knowledge distillation for anomaly detection. In *CVPR*, 2021. 2, 3, 6, 7
- [48] Thomas Schlegl, Philipp Seeböck, Sebastian M Waldstein, Georg Langs, and Ursula Schmidt-Erfurth. f-anogan: Fast unsupervised anomaly detection with generative adversarial networks. *Med Image Anal*, 2019. 3
- [49] Thomas Schlegl, Philipp Seeböck, Sebastian M Waldstein, Ursula Schmidt-Erfurth, and Georg Langs. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *IPMI*, 2017. 3
- [50] Hannah M Schlüter, Jeremy Tan, Benjamin Hou, and Bernhard Kainz. Natural synthetic anomalies for self-supervised anomaly detection and localization. In *ECCV*, 2022. 2
- [51] Bernhard Schölkopf, John C Platt, John Shawe-Taylor, Alex J Smola, and Robert C Williamson. Estimating the support of a high-dimensional distribution. *Neural Comput*, 2001. 3
- [52] Philipp Seeböck, Sebastian Waldstein, Sophie Klimscha, Bianca S Gerendas, René Donner, Thomas Schlegl, Ursula Schmidt-Erfurth, and Georg Langs. Identifying and categorizing anomalies in retinal imaging data. *arXiv preprint arXiv:1612.00686*, 2016. 1
- [53] Xian Tao, Xinyi Gong, Xin Zhang, Shaohua Yan, and Chandranath Adak. Deep learning for unsupervised anomaly localization in industrial images: A survey. *TIM*, 2022. 1, 2, 3, 6

- [54] Rui Tian, Zuxuan Wu, Qi Dai, Han Hu, Yu Qiao, and Yu-Gang Jiang. Resformer: Scaling vits with multi-resolution training. In *CVPR*, 2023. 4
- [55] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017. 2, 4
- [56] Guodong Wang, Shumin Han, Errui Ding, and Di Huang. Student-teacher feature pyramid matching for unsupervised anomaly detection. *arXiv preprint arXiv:2103.04257*, 2021. 3
- [57] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual recognition. *TPAMI*, 2020. 3
- [58] Junke Wang, Zuxuan Wu, Wenhao Ouyang, Xintong Han, Jingjing Chen, Yu-Gang Jiang, and Ser-Nam Li. M2tr: Multimodal multi-scale transformers for deepfake detection. In *ICMR*, 2022. 2, 4
- [59] Rui Wang, Dongdong Chen, Zuxuan Wu, Yinpeng Chen, Xiyang Dai, Mengchen Liu, Yu-Gang Jiang, Luwei Zhou, and Lu Yuan. Bevt: Bert pretraining of video transformers. In *CVPR*, 2022. 2
- [60] Zejia Weng, Xitong Yang, Ang Li, Zuxuan Wu, and Yu-Gang Jiang. Semi-supervised vision transformers. In *ECCV*, 2022. 4
- [61] Matthias Wieler and Tobias Hahn. Weakly supervised learning for industrial optical inspection. 2007. 1, 5, 8, 9, 10, 11
- [62] Zhen Xing, Qi Dai, Han Hu, Jingjing Chen, Zuxuan Wu, and Yu-Gang Jiang. Svformer: Semi-supervised video transformer for action recognition. In *CVPR*, 2023. 4
- [63] Minghui Yang, Peng Wu, Jing Liu, and Hui Feng. Memseg: A semi-supervised method for image surface defect detection using differences and commonalities. *arXiv preprint arXiv:2205.00908*, 2022. 2, 5, 9
- [64] Jihun Yi and Sungroh Yoon. Patch svdd: Patch-level svdd for anomaly detection and segmentation. In *ACCV*, 2020. 3
- [65] Sanyapong Youkachen, Miti Ruchanurucks, Teera Phatrapomnant, and Hirohiko Kaneko. Defect segmentation of hot-rolled steel strip surface by using convolutional auto-encoder and conventional image processing. In *IC-ICTES*, 2019. 2
- [66] Jongmin Yu, Du Yong Kim, Younkwan Lee, and Moongu Jeon. Unsupervised pixel-level road defect detection via adversarial image-to-frequency transform. In *IV*, 2020. 3
- [67] Jiawei Yu, Ye Zheng, Xiang Wang, Wei Li, Yushuang Wu, Rui Zhao, and Liwei Wu. Fastflow: Unsupervised anomaly detection and localization via 2d normalizing flows. *arXiv preprint arXiv:2111.07677*, 2021. 2, 3
- [68] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. Draem-a discriminatively trained reconstruction embedding for surface anomaly detection. In *ICCV*, 2021. 2, 5, 6, 7, 9, 10
- [69] Jianpeng Zhang, Yutong Xie, Guansong Pang, Zhibin Liao, Johan Verjans, Wenxing Li, Zongji Sun, Jian He, Yi Li, Chunhua Shen, et al. Viral pneumonia screening on chest x-rays using confidence-aware anomaly detection. *TMI*, 2020. 3
- [70] Ye Zheng, Xiang Wang, Rui Deng, Tianpeng Bao, Rui Zhao, and Liwei Wu. Focus your distribution: Coarse-to-fine non-contrastive learning for anomaly detection and localization. In *ICME*, 2022. 3