# Self-Supervised Masking for Unsupervised Anomaly Detection and Localization

Chaoqin Huang, *Student Member, IEEE*, Qinwei Xu, *Student Member, IEEE*,
Yanfeng Wang, *Member, IEEE*, Yu Wang, *Member, IEEE*, and Ya Zhang, *Member, IEEE*

*Abstract*—Recently, anomaly detection and localization in multimedia data have received significant attention among the machine learning community. In real-world applications such as medical diagnosis and industrial defect detection, anomalies only present in a fraction of the images. To extend the reconstruction-based anomaly detection architecture to the localized anomalies, we propose a self-supervised learning approach through *random masking* and then *restoring*, named <u>Self-Supervised Masking</u> (SSM) for unsupervised anomaly detection and localization. SSM not only enhances the training of the inpainting network but also leads to great improvement in the efficiency of mask prediction at inference. Through random masking, each image is augmented into a diverse set of training triplets, thus enabling the autoencoder to learn to reconstruct with masks of various sizes and shapes during training. To improve the efficiency and effectiveness of anomaly detection and localization at inference, we propose a novel progressive mask refinement approach that progressively uncovers the normal regions and finally locates the anomalous regions. The proposed SSM method outperforms several state-of-the-arts for both anomaly detection and anomaly localization, achieving 98.3% AUC on Retinal-OCT and 93.9% AUC on MVTec AD, respectively.

*Index Terms*—Anomaly detection, anomaly localization, self-supervised learning, progressive mask refinement.

## I. INTRODUCTION

ANOMALY detection and localization in multimedia data have received significant attention among the machine learning community, with broad application in medical diagnosis [1]–[4], defect detection in the factories [5], credit card fraud detection [6], and autonomous driving [7]. For most of the above applications, anomalous samples are remarkably scarce in the population. It is often prohibitive to collect a representative set of anomalous samples. As a result, many studies [8]–[10] have resorted to learning in the unsupervised setting, *i.e.,* training with normal samples only.

Along this line, previous studies attempt to first model the normal distribution through either one-class classification-based approaches [11]–[13], reconstruction-based approaches [14]–[16], or self-supervision-based approaches [17]–[19], and then detect the anomalies by identifying samples with different distributions than the models. With the recent advances in deep neural networks, reconstruction-based approaches have received increasing attention and shown great promise for unsupervised anomaly detection. The network, trained with normal data only, is assumed not generalizable to abnormal samples, and thus leads to high reconstruction errors for anomalous samples.

In real-world applications such as medical diagnosis [20] and industrial defect detection [21], anomalies only present in a fraction of the images. To extend the reconstruction-based architecture to the localized anomalies, a few recent methods [22], [23] leverage image inpainting to locate anomalies from their surrounding context, with the assumption that the difference between the masked region and its corresponding restoration is significant for anomalous regions. Inpainting here needs to solve the dual tasks of *masking* possible anomalies and *restoring* the masked regions. While accurate restoration is critical to finding the masks for anomalies, it also depends on the precise masking of anomalies. To break the circular dependency between masking and restoring, SCADN [22] introduces a fixed set of striped masks in multiple scales, and SMAI [23] leverages super-pixel segmentation to generate candidate masks, so that they can focus on fitting the restoration networks. Despite their promising results for anomaly localization, both methods require brutally traversing all possible masks, which is prohibitively time-consuming for real-world applications.

This paper explores a self-supervised learning approach through *random masking* and then *restoring*, named <u>Self-Supervised Masking</u> (SSM) hereafter. The key difference between SSM and the above inpainting-based approaches [22], [23] mainly lies in the way of masking, which not only enhances the training of the inpainting network but also greatly improves the efficiency of mask prediction at test time. **(i)** During each training epoch, a random mask is generated on-the-fly for each image, and the masked image is then fed to a conditional autoencoder with two prediction heads, one for image reconstruction and the other for mask reconstruction (Figure 1 (a)). By use of random masking, each image is augmented into a diverse set of training triplets <*masked image, mask, image*>, thus enabling the autoencoder to learn to reconstruct with masks of various sizes and shapes. **(ii)** During the inference, as it is impossible to brutally traverse all possible masks (*i.e.,* $2^n$ possible masks for an image of $n$ pixels), a novel *progressive mask refinement* approach is introduced to improve inference efficiency (Figure 1 (b)). SSM starts with a pair of complementary checkered masks, which jointly determine an initial mask according to the reconstruction results. The masks are iteratively refined and shrunk to the likely anomaly regions based on the reconstruction errors. Figure 2
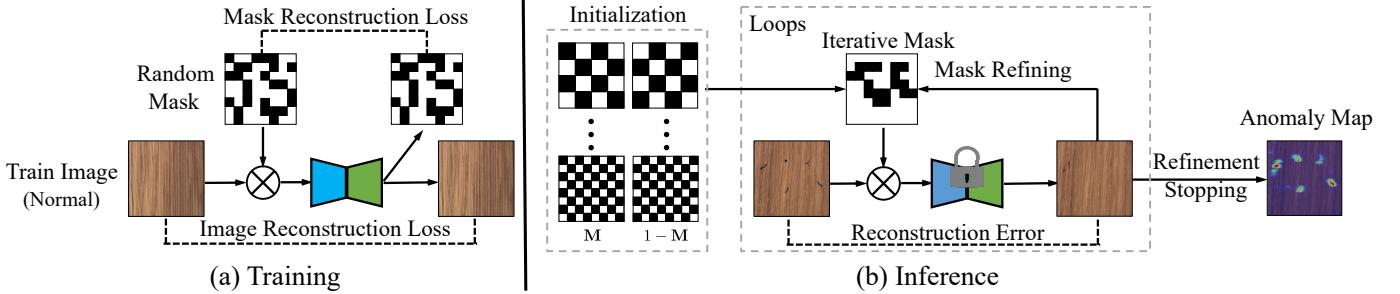
(a) Training | (b) Inference

Fig. 1: The overview of SSM. (a) During training, SSM leverages a conditional autoencoder to reconstruct the training images under randomly generated masks for only normal samples and also reconstruct the randomly generated masks, under a self-supervised learning paradigm. (b) During inference, SSM locates anomalies with a progressive mask refinement approach. An iterative mask is refined according to the feedback of the anomaly scores provided by the conditional autoencoder. SSM progressively uncovers the normal regions and finally locates the anomalous regions.
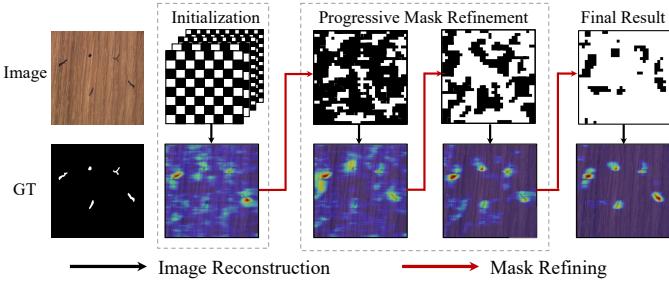


Fig. 2: Results of anomaly localization during the inference stage with SSM on the MVTec AD dataset. During the process of progressive mask refinement, SSM continually narrows the scopes of both masks (top) and localization maps (bottom) to the anomalous regions.

provides an illustration of the progressive anomaly localization process. To deal with anomalies of various shapes and sizes, the progressive mask refinement process is performed with initial masks of multiple scales and their ensemble results are used to detect the anomaly.

To validate the effectiveness of SSM, we experiment with two popular benchmark datasets, Retinal-OCT [20] for medical diagnosis and MVTec AD [21] for industrial defect localization. The experimental results have shown that SSM outperforms a number of state-of-the-art methods for both anomaly detection and localization, achieving 98.3% anomaly detection AUC on Retinal-OCT [20] and 93.9% anomaly localization AUC on MVTec AD [21], respectively.

The main contributions of the paper are summarized as follows:

- We propose a novel masking & restoring framework for unsupervised anomaly detection and localization named SSM. Through random masking, SSM enables the autoencoder to reconstruct the data with masks of various sizes and shapes, thus leading to a more powerful representation learning.
- To further improve the efficiency and effectiveness of anomaly detection and localization at inference, we propose a novel progressive mask refinement approach that progressively uncovers the normal regions and finally locates the anomalous regions.

## II. RELATED WORKS

### A. Unsupervised Anomaly Detection

Anomaly detection can be roughly divided into two classes: anomalous human behavior detection in videos [8], [24]–[32] and anomaly detection in still images (or outlier data detection). In this paper, we focus on anomaly detection and localization in images, especially for medical diagnosis [20] and industrial defect detection [21]. Compared with the supervised approaches, unsupervised anomaly detection and localization is to train with normal samples only, without any anomalous data, and no image-level annotation or pixel-level annotation is provided. Since no auxiliary information for anomalies is provided, approaches like zero-shot object detection [33], [34] are also infeasible for unsupervised anomaly detection. Under the unsupervised setting, the majority of the research in image anomaly detection can be broadly categorized as one-class classification-based approaches, reconstruction-based approaches, generative adversarial network (GAN)-based approaches, and self-supervision-based approaches.

*1) One-class Classification-based Approaches:* One-class classification is referred to as the problem of learning a description of a set of data instances to detect whether new instances conform to the training data or not. It assumes that all normal instances can be summarized by a compact model, to which anomalies do not conform [35]–[38]. In OC-SVM [11], the normal samples are mapped to the high-dimensional feature space through a kernel function to get better aggregated. It learns a hyperplane that maximizes a margin between training data instances and the origin. To better aggregate the mapped data in latent space, Deep SVDD [13] optimizes the neural network by minimizing the volume of a hyper-sphere that encloses the network representations of the data. However, the one-class models may fail for datasets with complex distributions within the normal class.

*2) Reconstruction-based Approaches:* Reconstruction-based anomaly detection approaches [14], [15], [39]–[41] aim to learn the low-dimensional feature representation space on which the given normal data instances can be well reconstructed. The heuristic for using this technique in anomaly detection is that the learned feature representations are trained to learn regularities of the data. From this representation, anomalies are

difficult to be reconstructed and thus have large reconstruction errors. DAE [42] firstly applies the autoencoder to anomaly detection. To improve DAE, Nicolau *et al.* [43] introduce a density estimator, Kernel Density Estimation (KDE) [44], to model the density from the hidden layer of autoencoders. By placing a threshold on the density of the normal data, query points below the threshold are classed as anomalies.

Recently, a deep autoencoder is adopted to improve the feature representation abilities [45]–[48]. In the same vein, MemAE [16] augments the autoencoder with a memory module to highlight reconstructed errors on anomalies. To capture the information of image texture and structure, P-Net [47] proposes to leverage the relation between the image texture and structure to enlarge the structure difference, which can also be used as a metric for normality measurement. Based on [47], MemSTC [48] proposes a structure-texture correspondence memory module to reconstruct image texture from its structure, where a memory mechanism is used to characterize the mapping from the normal structure to its corresponding normal texture.

To detect anomalies that present in a fraction of an image, recent methods [22], [23] leverage image inpainting to extend the reconstruction-based approach to locate anomalies from their surrounding context. The difference between the masked region and its corresponding restoration is assumed to be significant for anomalous regions. SCADN [22] leverages a fixed set of striped masks in multiple scales and traverses all possible masks at inference. Similarly, SMAI [23] leverages a superpixel masking and inpainting framework to identify and locate anomalies, where an inpainting module is trained to learn the spatial and texture information of the normal samples through random superpixel masking and restoration. Although the superpixel masks in SMAI may be more appropriate than the striped masks in SCADN, it takes longer for inference due to the need of traversing more possible masks (77 possible masks in total). This paper introduces a progressive mask refinement approach to avoid brutally traversing all possible masks, which improves inference efficiency.

*3) Generative Adversarial Network-based Approaches:* By assuming that normal data instances can be better generated than anomalies from the latent feature space of the generative network, adversarial training is employed [8], [30], [49]. These approaches generally aim to learn a latent feature space of a generative network so that the latent space well captures the normality underlying the given data. Along this line, to improve the generator's robustness against noises, ALOCC adds Gaussian noises to the inputs to form the training normal samples [8]. GANomaly [50] leverages another encoder to embed the generated results to a subspace, and calculates the anomaly scores in the subspace but not in the image space like ALOCC. OCGAN [9] further applies two adversarial discriminators and a classifier on a denoising autoencoder. By adding constraints and forcing latent codes to reconstruct examples like the normal data, anomalies show higher reconstruction errors. Similarly, adVAE [51] employs adversarial training within a variational autoencoder framework under the assumption that normal and anomalous data follows different Gaussian distributions. However, the generator networks can be misled and thus generate data instances out of the manifold of normal

instances, especially when the distribution of the given dataset is complex or the training data contains unexpected outliers.

*4) Self-supervision-based Approaches:* Recently, self-supervised learning has been widely used as it benefits many downstream tasks like classification [52], [53], detection [34], [54], segmentation [53], and tracking [55]. Among various self-supervised tasks, image restoration or inpainting has been widely adopted, where the network is forced to learn rich and robust feature embeddings. Pathak *et al.* [56] propose to remove arbitrary shapes from the input images. Those shapes are obtained as objects in the PASCAL VOC 2012 dataset [57] and pasted in arbitrary places in the other images. In [58], a low resolution but intact version of the original image is further fed to the network to guide the restoration. Based on [56], [59] further proposes to remove and restore the internal representations and designs a feature attention module to improve the feature robustness.

Self-supervision-based anomaly detection approaches [17], [19], [60] learn the representations of the normal data under a self-supervised learning paradigm with different self-supervisions [52]–[54]. Models are optimized with different surrogate tasks. Then anomalies can be separated under the assumption that anomalies will result differently in the corresponding surrogate task. Similarly, GeoTrans applies dozens of image geometric transforms and creates a self-labeled dataset for transformation classification [17], [53]. It assumes that transformations applied on anomalous data can not be classified properly. Wang *et al.* [18] apply Jigsaw puzzles [52] to extend the above self-labeled dataset. [61] learns representations by classifying data between different types of CutPaste, a set of data augmentations, and then utilizes a sliding window for anomaly localization although it is time-consuming. ARNet [19] applies image restoration as the self-supervision, assuming that the model is able to learn semantic features during the restoration process. GP [62] proposes a patch-based approach that considers both the global and local information and utilizes the discrepancy between global and local features as the anomaly score. Based on the knowledge distillation approach, US [60] and MKD [63] train the student networks to regress the output of a descriptive teacher network that was pre-trained on a large dataset. Anomalies are detected when the outputs of the student networks differ from that of the teacher network. This happens when they fail to generalize outside the manifold of anomaly-free training data. In this paper, self-supervised learning is applied in the proposed approach for both conditional image reconstruction and mask reconstruction.

### B. Progressive Refinement

The progressive refinement network has been explored in many supervised tasks, such as supervised image matting [64], person re-identification [65], and temporal action detection [66], motivated by the thought of progressive learning [67], [68]. For example, PBRNet [66] is equipped with three cascaded detection modules for progressive localizing action boundaries more and more precisely. MGMatting [64] proposes a progressive refinement network for image matting, which encourages the matting model to provide self-guidance to
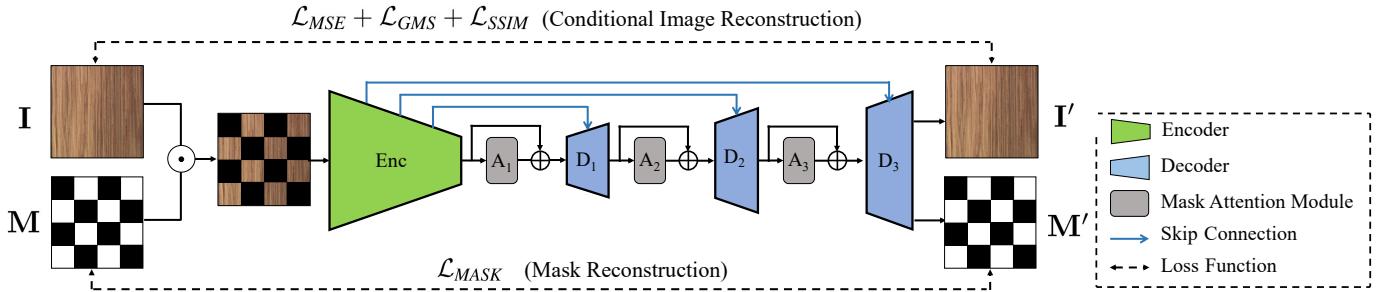
Fig. 3: Model architecture of the conditional autoencoder for the self-supervised masking training. An autoencoder $\{\text{Enc}, D_1, D_2, D_3\}$ takes $\mathbf{I} \odot \mathbf{M}$ as input and outputs both reconstructed images $\mathbf{I}'$ (top) and the reconstructed masks $\mathbf{M}'$ (bottom). Some mask attention modules $\{A_1, A_2, A_3\}$ distributed across the layers of the decoder (see Figure 4 for more details).
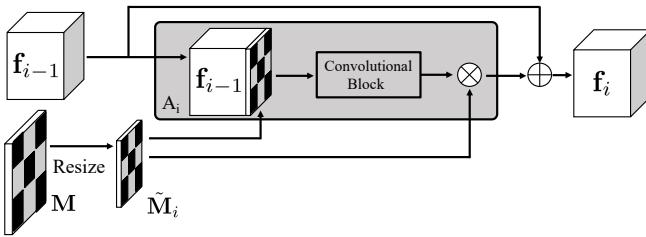


Fig. 4: Mask attention module (MAM) of conditional autoencoder for self-supervised masking training. The mask attention module $A_i$ takes the feature $\mathbf{f}_{i-1}$ as input and output $\mathbf{f}_i$ for the next decoder part. Following [59], we use a residual block design and gate the output of the last convolutional layer with the resized mask $\tilde{\mathbf{M}}_i$.

progressively refine the uncertain regions through the decoding process in multiple layers of the feature hierarchy. In this paper, we leverage progressive mask refinement for unsupervised anomaly detection. Different from the multiple feature hierarchy refinement for supervised tasks, the proposed progressive mask refinement is a procedure at inference, which takes the continuously refined mask as the input and reuses the trained conditional autoencoder for image reconstruction.

## III. METHOD

The key novelties of SSM include *random masking and restoring* at training and *progressive mask refinement* at inference. As the training data are all normal, a random mask is generated to formulate the inpainting task for each image. A conditional autoencoder is then leveraged to restore the masked region. At inference, to avoid brutally traversing all possible masks, a progressive mask refinement approach is proposed to improve the efficiency of anomaly detection and localization.

### A. Random Masking

To define the inpainting area of the training images, each input image is decomposed into $\frac{H}{k} \times \frac{W}{k}$ grids, where $H$ and $W$ are the height and width of the images, and $k$ here controls the granularity of the grid. Each grid consists of a square of $k \times k$ pixels and is set as the basic unit of the masks. The match of the grid size $k$ and the size of the anomaly is expected to significantly influence the performance of anomaly

detection algorithms. If $k$ is much larger or much smaller than the anomaly, it is generally infeasible to obtain accurate reconstruction. Since anomalies could come in various sizes, and there is no way to know the size of the anomalies as prior, we here consider detecting with multiple scales by varying the values of $k$. Particularly, the size $k$ is sampled from a set $K = \{k_i\}_{i=1:N_k}$, where $N_k$ is the set cardinality. In our implementation, we use $K = \{4, 8, 16, 32\}$ as it covers a wide range of anomaly scales.

To enlarge the exploration space of the masks, during each training epoch, a random mask is generated on-the-fly for each image. Each grid is then randomly chosen to be masked or to be kept, and the resulting mask matrix is denoted as $\mathbf{M}$. In this way, a diverse set of random masks with various sizes and shapes are generated. Through such random masking, each image is augmented into a diverse set of training triplets $< \tilde{\mathbf{I}}, \mathbf{M}, \mathbf{I} >$, where $\mathbf{I}$ is the input image, $\mathbf{M}$ is the generated spatial mask, $\tilde{\mathbf{I}} = \mathbf{I} \odot \mathbf{M}$ is the resulting masked image, and $\odot$ is the element-wise product in the spatial domain (the mask is replicated along the channel dimension), thus enabling to learn to reconstruct with masks of various sizes and shapes.

### B. Restoration Network

The overall architecture of the restoration network is shown in Figure 3. The backbone of the restoration network is a conditional autoencoder. Different from the vanilla autoencoder, the conditional autoencoder is used to encode unmasked regions and fill in masked regions under a certain mask matrix [69]. The condition lies in that the proposed image reconstruction network is mask-guided but not simply reconstructing the whole image. Different from vanilla image reconstruction-based anomaly detection methods, we assume that the discrepancy between the masked region and its corresponding restoration is significant for detecting anomalies.

For better representation learning, we further add a mask reconstruction branch in our network so that two prediction heads are associated with the conditional autoencoder, one for image reconstruction and the other for mask reconstruction:

$$\mathbf{I}', \mathbf{M}' = \text{Dec}[\text{Enc}(\mathbf{I} \odot \mathbf{M})], \tag{1}$$

where Enc is the encoder, Dec is the decoder, $\mathbf{M}'$ is the reconstructed mask and $\mathbf{I}'$ is the reconstructed image. Skip-connections between the encoder and decoder are added to

facilitate the backpropagation of gradients and improve the performance of image reconstruction.

Inspired by [59], to improve the robustness and reconstruction ability of the model in the manner of self-supervised learning, we further add a mask attention module (MAM) to the conditional autoencoder. The architecture of MAM is shown in Figure 4. The decoder Dec is split into three sub-networks $\{D_1, D_2, D_3\}$, and the mask attention module, $A = \{A_1, A_2, A_3\}$, is added in front of each of the sub-networks for the decoder network, as shown in Figure 3. For the $i$-th mask attention module $A_i$, we down-sample the mask $\mathbf{M}$ to $\tilde{\mathbf{M}}_i$ with the nearest neighbor method and match the spatial dimension of the corresponding input feature map $\mathbf{f}_{i-1}$. The output $\mathbf{f}_i$ of the $i$-th mask attention module $A_i$ is:

$$\mathbf{f}_i = \mathbf{f}_{i-1} + \phi(\mathbf{C}(\mathbf{f}_{i-1}, \tilde{\mathbf{M}}_i)) \odot \tilde{\mathbf{M}}_i, \quad (2)$$

where $\mathbf{C}(\cdot, \cdot)$ is a concatenation layer and $\phi(\cdot)$ is a convolutional block following [59], and $\odot$ denotes the element-wise product in the spatial domain (the mask is replicated along the channels). With the mask attention module, the model pays more attention to learn the context feature, thus the model's image restoration ability can be significantly improved.

### C. Loss Functions

Given the input image $\mathbf{I}$, and the reconstructed image $\mathbf{I}'$, as the network only targets to restore the masked regions, the unmasked regions are copied from the original images with an identity function:

$$\hat{\mathbf{I}} = \mathbf{I}' \odot (1 - \mathbf{M}) + \mathbf{I} \odot \mathbf{M}. \quad (3)$$

For image reconstruction, we aim to measure the difference between $\mathbf{I}$ and $\hat{\mathbf{I}}$ and consider the following set of loss functions.
**Mean Square Error:** The mean square error (MSE) is typically used for training an autoencoder.

$$\mathcal{L}_{MSE} = \|\mathbf{I} - \hat{\mathbf{I}}\|_2^2. \quad (4)$$

However, this loss assumes the independence between neighboring pixels, which may be incorrect in some situations. To solve this problem, this paper further introduces several losses that penalize structural differences between the reconstructed regions and the input regions, *i.e.,* a gradient magnitude similarity (GMS) loss [70] and a structured similarity index (SSIM) loss [71]. SSIM and GMS are both patch similarity metrics that focus on different image properties.
**Gradient Magnitude Similarity Loss:** The gradient difference map is defined as:

$$\mathcal{L}_{GMS} = \frac{1}{H \times W} \sum_{i=1}^{H} \sum_{j=1}^{W} \mathbf{1} - GMS(\mathbf{I}, \hat{\mathbf{I}})_{(i,j)}, \quad (5)$$

$$GMS(\mathbf{I}, \hat{\mathbf{I}}) = \frac{1}{3} \sum_{c=1}^{3} GMS_c(\mathbf{I}^c, \hat{\mathbf{I}}^c) \in \mathbb{R}^{H \times W}, \quad (6)$$

where $\mathbf{1}$ is a matrix of ones. $\mathbf{I}^c$ and $\hat{\mathbf{I}}^c$ are $c$-th color channel of the original image and the reconstructed image, respectively.

$GMS_c(\mathbf{I}^c, \hat{\mathbf{I}}^c) \in \mathbb{R}^{H \times W}$ is the gradient magnitude similarity map for color channel $c$:

$$GMS_c(\mathbf{I}^c, \hat{\mathbf{I}}^c) = \frac{2g(\mathbf{I}^c)g(\hat{\mathbf{I}}^c) + a}{g(\mathbf{I}^c)^2 + g(\hat{\mathbf{I}}^c)^2 + a}, \quad (7)$$

$$g(\mathbf{I}^c) = \sqrt{(\mathbf{I}^c * \mathbf{h}_x)^2 + (\mathbf{I}^c * \mathbf{h}_y)^2}, \quad (8)$$

where $a$ is a constant ensuring numerical stability. $\mathbf{h}_x$ and $\mathbf{h}_y$ are $3 \times 3$ Prewitt filters along the $x$ and $y$ dimensions and $*$ is the convolution operation. The GMS loss is differentiable and has been widely used by image inpainting and image super-resolution tasks [72], [73].
**Structured Similarity Index (SSIM):** The SSIM loss is defined as:

$$\mathcal{L}_{SSIM}(\mathbf{I}, \hat{\mathbf{I}}) = \frac{1}{H \times W} \sum_{i=1}^{H} \sum_{j=1}^{W} \mathbf{1} - SSIM(\mathbf{I}, \hat{\mathbf{I}})_{(i,j)}, \quad (9)$$

where $SSIM(\mathbf{I}, \hat{\mathbf{I}})_{(i,j)}$ is the SSIM [71] value between two patches of $\mathbf{I}$ and $\hat{\mathbf{I}}$ centered at $(i, j)$.

The mask reconstruction simply adopts the $L_2$ loss:

$$\mathcal{L}_{MASK} = \|\mathbf{M} - \mathbf{M}'\|_2^2. \quad (10)$$

Finally, the total loss for training SSM is formulated as:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{MSE} + \lambda_2 \mathcal{L}_{GMS} + \lambda_3 \mathcal{L}_{SSIM} + \lambda_4 \mathcal{L}_{MASK}, \quad (11)$$

where $\lambda_1$, $\lambda_2$, $\lambda_3$, and $\lambda_4$ are the individual loss weights.

### D. Progressive Mask Refinement at Inference

At inference time, one additional challenge emerges, *i.e.,* how to set up the mask for conditional reconstruction? It is desired that the mask covers only the anomalous regions but not the normal region, which is in fact a dilemma: if we have a perfect mask, we can correctly locate the anomalies; otherwise, the reconstruction result may not correctly manifest the anomalies. To break the dilemma, existing inpainting-based approaches [22], [23] simply brutally traverse all possible masks, but at the cost of limiting the search space for masks. In our case, with random masking, it is impossible to brutally traverse all possible masks. A novel progressive mask refinement approach, consisting of two stages, *i.e.*, mask initialization and mask refinement, is introduced to improve the efficiency for inference.

Given the input image $\mathbf{I}$ and the reconstructed image $\hat{\mathbf{I}}$, an error function $f(\cdot, \cdot)$ is introduced as follows.

$$f(\mathbf{I}, \hat{\mathbf{I}}) = L_2(\mathbf{I}, \hat{\mathbf{I}}) + (1 - GMS(\mathbf{I}, \hat{\mathbf{I}})) + (1 - SSIM(\mathbf{I}, \hat{\mathbf{I}})), \quad (12)$$

where $L_2(\cdot, \cdot)$, $GMS(\cdot, \cdot)$, $SSIM(\cdot, \cdot)$ are the per-pixel $L_2$, GMS and SSIM score maps, respectively. This error function is then leveraged for both anomaly score calculation and mask refinement. By calculating the per-pixel-based error scores and treating regions with large scores as potential anomalous, the masks are iteratively refined and gradually shrunk to the possible anomaly regions.

(a) $4 \times 4$ patches     (b) $8 \times 8$ patches

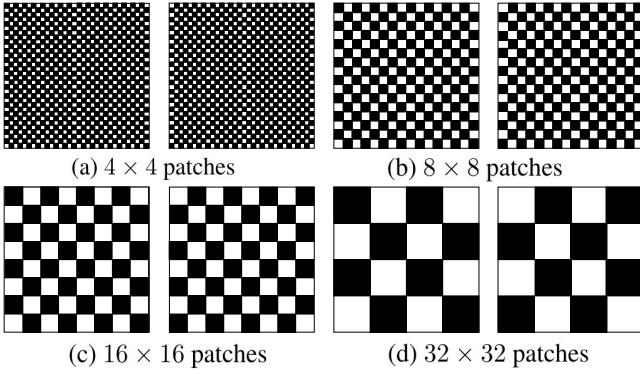(c) $16 \times 16$ patches     (d) $32 \times 32$ patches

Fig. 5: Mask initialization during the inference. These initialized masks are complementary with each other for each size of patches and thus cover all pixels in the image, which avoids missing possible anomalous areas.

*1) Mask Initialization:* Given a test image, assuming no prior information about the anomalous region at the beginning, as initialization, we start with a set of multi-scaled masks $\mathcal{M}^0$, which consists of eight checkerboard-like matrices in different scales. Figure 5 shows the set of initialization masks used in our experiments, where the grid size $k \in K, K = \{4, 8, 16, 32\}$ in our experiments. For each grid size $k$, the initialization masks contain a pair of complementary masks which jointly cover all pixels in the image, thus avoiding missing any possible anomalous areas.

For each initialization mask, the model reconstructs the pixels in masked regions of the test images based on the corresponding pixels in un-masked regions. The anomaly score map is obtained for each reconstructed image using the error function $f(\cdot, \cdot)$ as defined in Eq. 12. The score maps for different masks are average into a single one to obtain the initialized anomaly score map $\mathbf{S}^0$.

*2) Mask Refinement:* The purpose of mask refinement is to remove the masked areas likely corresponding to the normal regions, so that the conditional image reconstruction network pays more attention to the remaining anomalous regions. At each iteration, the reconstruction error map is leveraged to refine the mask, by considering regions with smaller errors as normal and removing them from the mask for the next iteration.

When most of the regions covered by the mask are anomalous regions, providing more image information can not reduce the reconstruction error of anomalous regions significantly, and the corresponding mask is unchanged/converged. At this point, we terminate the inference stage, and obtain the final masks. Finally, when this approach ends, the mask is expected to cover only the anomalous parts of the image.

The corresponding algorithm of the mask refinement is shown in Algorithm 1. The mask is updated in the unit of patches to make the algorithm more stable and reduce the number of iterations. Thus, given the grid size $k \in K$, the image $\mathbf{I}$ is split into $N_k$ $k \times k$ patches $p_1, p_2, \cdots, p_{N_k}$. Then, for each patch $p$, we calculate the average reconstruction error $\epsilon_p$ on $p$ and use $\epsilon_p$ as a credential to update the mask $\mathbf{M}$ according to a threshold $\eta$. The threshold $\eta$ is set to be the

---

**Algorithm 1:** Progressive Mask Refinement

**Input:** input $\mathbf{I} \in \mathbb{R}^{H \times W \times C}$,
    conditional autoencoder $\text{Dec}[\text{Enc}(\cdot)]$,
    anomaly score map calculation function $f(\cdot, \cdot)$,
    an initialized score map $\mathbf{S}^0$,
    patch size $k$, a threshold $\eta$.
**Output:** anomaly score $\epsilon_k$ (for detection),
    anomaly score map $\mathbf{S}^k$ (for localization).

1   $\mathbf{S}^k = \mathbf{S}^0$
2   $N_k = \frac{H}{k} \times \frac{W}{k}$         # numbers of the patches
3   split image $\mathbf{I}$ into $N_k$ $k \times k$ patches $p_1, p_2, \cdots, p_{N_k}$
4   **repeat**
5     # mask refinement
6     **for** *each patch $p$* **do**
7       $\epsilon_p = \frac{1}{k^2} \sum_{(x,y) \in p} \mathbf{S}^k_{(x,y)}$
8       **for** *each pixel $(x, y)$ in $p$* **do**
9         $\mathbf{M}_{(x,y)} = \begin{cases} 0 & \text{if } \epsilon_p > \eta \\ 1 & otherwise \end{cases}$
10     # score calculation
11     $\mathbf{I}' = \text{Dec}[\text{Enc}(\mathbf{I} \odot \mathbf{M})]$
12     $\hat{\mathbf{I}} = \mathbf{I}' \odot (1 - \mathbf{M}) + \mathbf{I} \odot \mathbf{M}$
13     $\mathbf{S}^k = f(\mathbf{I}, \hat{\mathbf{I}})$       # anomaly score map
14   **until** $\mathbf{M}$ *converges*;
15   $\epsilon_k = \frac{1}{\sum_{i,j} \mathbf{M}_{(i,j)}} \sum_{x,y} \mathbf{S}^k_{(x,y)}$     # anomaly score
16   **return** anomaly score $\epsilon_k$, anomaly score map $\mathbf{S}^k$

---

maximum error in the validation set. Since the validation set contains only normal samples, this threshold can be considered as a rough boundary between the normal and anomalous cases. We terminate this approach until $\mathbf{M}$ converges. For each grid size $k$, we obtain an anomaly score $\epsilon_k$ and an anomaly score map $\mathbf{S}^k$ according to the mask refinement algorithm.

To obtain the final anomaly score map $\mathbf{S}^{final}$, we average all the anomaly score maps $\mathbf{S}^k$ provided by those using different grid sizes of masks. Similarly, the anomaly score $\epsilon$ is computed as the average of $\epsilon_k$ under several grid sizes $k \in K$.

*E. Discussion*

In our framework, the progressive mask refinement approach is only deployed in the inference stage. A natural question is: Why not use it in the training stage? To answer this question, let us recall the definition of anomaly detection. In the training stage, with only normal data provided, it is meaningless to update the mask: the entire image is expected to be well recovered, ending up with an empty mask after the progressive mask refinement. This also violates the conditional reconstruction principle we proposed. Thus, a better choice is reconstructing the images with randomly generated masks. With a larger exploration space of masks, SSM leads to a more powerful representation learning. This design alleviates the impact of the lack of anomalous data during training. We believe that this is a feasible research avenue for unsupervised anomaly detection and localization.

TABLE I: Results of **anomaly detection** in terms of AUC in % on the Retinal-OCT dataset, comparing with several state-of-the-arts. The best-performing method is in bold.

| Method | GeoTrans [17] | EGBAD [49] | f-AnoGAN [74] | DSVDD [13] | Pix2Pix [75] | AE [45] | AnoGan [41] | EnGAN [76] | C-GAN [77] | V-GAN [1] | P-Net [47] | MKD [63] | SSM (ours) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| AUC (%) | 60.1 | 61.0 | 66.6 | 74.4 | 79.4 | 82.1 | 84.8 | 86.9 | 87.4 | 90.6 | 92.9 | 97.0 | **98.3** |

TABLE II: Results of **anomaly detection** on the MVTec AD dataset. Results are listed as AUC in % and are marked individually for each class. An average score over all classes is also reported in the last row. Results of GeoTrans, GANomaly and ARNet are borrowed from [19]. Results of OCGAN, ALOCC, DAE, MemAE and SCADN are borrowed from [22]. Results of MKD are borrowed from [63]. Results of MemSTC are borrowed from [48]. The best-performing method is in bold.

| Category | OCGAN [9] | ALOCC [8] | GeoTrans [17] | DAE [42] | MemAE [16] | GANomaly [50] | SCADN [22] | ARNet [19] | US [60] | MKD [63] | MemSTC [48] | SSM (ours) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Bottle | 59.2 | 46.0 | 74.4 | 86.0 | 93.0 | 89.2 | 95.7 | 94.1 | 99.0 | 99.4 | 97 | **99.9** |
| Cable | 49.6 | 53.1 | 78.3 | 64.8 | 78.5 | 74.5 | 85.6 | 83.2 | 86.2 | **89.2** | 81 | 77.3 |
| Capsule | 71.4 | 48.7 | 67.0 | 53.4 | 73.5 | 73.2 | 76.5 | 68.1 | 86.1 | 80.5 | 87 | **91.4** |
| Carpet | 34.8 | 42.3 | 43.7 | 58.8 | 38.6 | 69.9 | 50.4 | 70.6 | **91.6** | 79.3 | 61 | 76.3 |
| Grid | 85.5 | 78.1 | 61.9 | 85.8 | 80.5 | 70.8 | 98.3 | 88.3 | 81.0 | 78.0 | 99 | **100** |
| Hazelnut | 75.3 | **99.3** | 35.9 | 51.3 | 76.9 | 78.5 | 83.3 | 85.5 | 93.1 | 98.4 | 98 | 91.5 |
| Leather | 62.4 | 76.8 | 84.1 | 49.7 | 42.3 | 84.2 | 65.9 | 86.2 | 88.2 | 95.1 | 87 | **99.9** |
| Metal Nut | 29.5 | 70.5 | 81.3 | 79.3 | 65.4 | 70.0 | 62.4 | 66.7 | 82.0 | 73.6 | 82 | **88.7** |
| Pill | 70.2 | 72.6 | 63.0 | 69.3 | 71.7 | 74.3 | 81.4 | 78.6 | 87.9 | 82.7 | 87 | **89.1** |
| Screw | 50.5 | 99.5 | 50.0 | 71.9 | 25.7 | 74.6 | 83.1 | **100** | 54.9 | 83.3 | 99 | 85.0 |
| Tile | 80.6 | 52.6 | 41.7 | 89.4 | 71.8 | 79.4 | 79.2 | 73.5 | **99.1** | 91.6 | 98 | 94.4 |
| Toothbrush | 59.4 | 64.2 | 97.2 | 94.2 | 96.7 | 65.3 | 98.1 | **100** | 95.3 | 92.3 | **100** | **100** |
| Transistor | 47.7 | 75.1 | 86.9 | 37.6 | 79.1 | 79.2 | 86.3 | 84.3 | 81.8 | 85.6 | 89 | **91.0** |
| Wood | 95.9 | 27.9 | 61.1 | 88.2 | 95.4 | 83.4 | 96.8 | 92.3 | 97.7 | 94.3 | **98** | 95.9 |
| Zipper | 36.4 | 54.7 | 82.0 | 81.9 | 71.0 | 74.5 | 84.6 | 87.6 | 91.9 | 93.2 | 93 | **99.9** |
| Mean | 60.6 | 64.1 | 67.2 | 70.7 | 70.7 | 76.2 | 81.8 | 83.9 | 87.7 | 87.7 | 90 | **92.0** |

## IV. EXPERIMENTS

In this section, our method is applied to the Retinal-OCT dataset [20] for unsupervised anomaly detection on medical diagnosis, and the MVTec Anomaly Detection [21] dataset for both image-level anomaly detection and pixel-level anomaly localization, compared with state-of-the-art methods. To evaluate the effectiveness of our method, we further conduct ablation studies under different loss functions and architecture designs of SSM. Finally, we provide visualization analysis to illustrate the effectiveness of the proposed progressive mask refinement approach.

### A. Experimental Setups

*1) Evaluation Protocols:* We quantify the model performance using the area under the Receiver Operating Characteristic (ROC) curve metric (AUC). This evaluation protocol allows comparison using different thresholds on the anomaly score. It is commonly adopted as the performance measurement in anomaly detection tasks.

*2) Datasets:* We conduct experiments on two real-world anomaly detection datasets, which are related to the medical diagnosis and the industrial defect detection:

**Retinal-OCT dataset** [20] is a recent dataset for detecting abnormalities in retinal optical coherence tomography (OCT) images. It contains 84,495 high-resolution clinical images. It has three small disease classes (CNV, DME, DRUSEN) and a large class of disease-free images. We use the large disease-free class as normal data; then, we use the three disease classes together as a single anomalous class. The training-testing split is set the same as the original dataset.

**MVTec Anomaly Detection dataset** [21] comprises 15 categories with 3629 images for training and validation and 1725 images for testing. The training set contains only normal images without defects. The test set contains images containing various kinds of defects and defect-free images. In total, 73 different defect types are present, on average five per category. Five categories cover different types of regular (carpet, grid) or random (leather, tile, wood) textures, while the remaining ten categories represent various types of objects. All image resolutions are in the range between $700 \times 700$ and $1024 \times 1024$ pixels. Pixel-precise ground truth labels for each defective image region are provided. The dataset contains almost 1900 manually annotated regions.

*3) Baselines:* For anomaly detection, we consider several state-of-the-art methods as baselines. For the one-class classification-based method, we consider DSVDD [13] as the baseline. For the reconstruction-based mathods, DAE [42], AE [45], MemAE [16], SCADN [22] and MemSTC [48] are considered. For GAN-based methods, Pix2Pix [75], AnoGAN [41], C-GAN [77], V-GAN [1], EGBAD [49], ALOCC [8], GANomaly [50], f-AnoGAN [74], OCGAN [9] and EnGAN [76] are considered. For other self-supervision-based methods, we consider GeoTrans [17] and ARNet [19]. We also consider the knowledge distillation-based methods, US [60] and MKD [63], which use additional data to pre-train the networks, to establish strong baselines. Especially, in

TABLE III: Results of **anomaly localization** on the MVTec AD dataset. Results are listed as AUC in % and are marked individually for each class. An average score over all classes is also reported in the last row. Results of DAE, OCGAN, MemAE and SCADN are borrowed from [22]. Results of AnoGAN, CNN-dict and SMAI are borrowed from [23]. Results of GDR and MKD are borrowed from [63]. Results of GP are borrowed from [62]. The best-performing method is in bold.

| Category | DAE [42] | OCGAN [9] | MemAE [16] | AnoGAN [41] | SCADN [22] | CNN-dict [78] | SMAI [23] | GDR [79] | MemSTC [48] | MKD [63] | GP [62] | SSM (ours) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Bottle | 54.4 | 56.7 | 72.4 | 86 | 69.6 | 78 | 86 | 92.2 | 87.2 | **96.3** | 93 | 95.9 |
| Cable | 53.5 | 56.4 | 81.4 | 78 | 81.4 | 79 | 92 | 91.0 | 91.2 | 82.4 | **94** | 82.1 |
| Capsule | 54.2 | 63.7 | 67.3 | 84 | 68.7 | 84 | 93 | 91.7 | 91.2 | 95.9 | 90 | **98.4** |
| Carpet | 52.8 | 54.6 | 57.4 | 54 | 64.9 | 72 | 88 | 73.5 | 85.7 | 95.6 | **96** | 94.4 |
| Grid | 55.0 | 65.2 | 46.8 | 58 | 79.6 | 59 | 97 | 96.1 | 93.9 | 91.8 | 78 | **99.0** |
| Hazelnut | 66.4 | 84.1 | 84.6 | 87 | 88.4 | 72 | 97 | **97.6** | 96.1 | 94.6 | 84 | 97.4 |
| Leather | 78.3 | 74.9 | 68.6 | 64 | 76.3 | 87 | 86 | 92.5 | 95.7 | 98.1 | 90 | **99.6** |
| Metal Nut | 53.9 | 53.4 | 76.9 | 76 | 75.4 | 82 | **92** | 90.7 | 89.0 | 86.4 | 91 | 89.6 |
| Pill | 55.5 | 59.6 | 73.7 | 87 | 74.7 | 68 | 92 | 93.0 | 93.1 | 89.6 | 93 | **97.8** |
| Screw | 57.0 | 70.8 | 73.2 | 80 | 87.6 | 87 | 96 | 94.5 | 90.1 | 96.0 | 96 | **98.9** |
| Tile | 63.0 | 59.2 | 64.7 | 50 | 67.7 | **93** | 62 | 65.4 | 85.9 | 82.8 | 80 | 90.2 |
| Toothbrush | 61.6 | 76.3 | 88.6 | 90 | 90.1 | 77 | 96 | 98.5 | 95.2 | 96.1 | 96 | **98.9** |
| Transistor | 53.2 | 58.2 | 71.4 | 80 | 68.9 | 66 | 85 | 91.9 | 86.9 | 76.5 | **100** | 80.1 |
| Wood | 61.2 | 65.5 | 65.2 | 62 | 67.2 | **91** | 80 | 83.8 | 85.1 | 84.8 | 81 | 86.9 |
| Zipper | 53.6 | 62.4 | 64.3 | 78 | 67.0 | 76 | 90 | 86.9 | 89.4 | 93.9 | **99** | 99.0 |
| Mean | 58.2 | 64.1 | 70.4 | 74 | 75.2 | 78 | 89 | 89.3 | 90.4 | 90.7 | 91 | **93.9** |

medical diagnosis, we consider P-Net [47] as a baseline, which uses additional medical domain knowledge and uses a large amount of additional data and annotations for the training.

For anomaly localization, we follow [23], [62], [63] and consider DAE [42], OCGAN [9], MemAE [16], AnoGAN [41], SCADN [22], CNN-dict [78], SMAI [23], GDR [79], MKD [63] and GP [62] as baselines.

*4) Model Configuration:* For the encoder-decoder structure for SSM, we follow the settings in [19], [23], [47], [80] and add skip-connections between some layers in the encoder and corresponding decoder layers to facilitate the backpropagation of the gradient in an attempt to improve the performance of image reconstruction. The individual loss weights $\lambda_1$, $\lambda_2$, $\lambda_3$ and $\lambda_4$ are set to 1 as default. We use Adam [81] optimizer with the weight decay of $1e^{-5}$. Other hyperparameters are default in Pytorch. SSM is trained using a batch size of 8 for 300 epochs with one NVIDIA GTX 3090. The learning rate is initially set to $1e^{-4}$, and is divided by 2 every 50 epochs.

### B. Comparison with State-of-the-art Methods

In this section, we show quantitative results of the proposed SSM, comparing with several state-of-the-art methods on the Retinal-OCT dataset and MVTec AD dataset.

**Results on Retinal-OCT dataset.** Table I shows the AUC results of anomaly detection on Retinal-OCT dataset, comparing with several state-of-the-arts, including GeoTrans [17], EGBAD [49], f-AnoGAN [74], DSVDD [13], Pix2Pix [75], AE [45], AnoGAN [41], EnGAN [76], C-GAN [77], V-GAN [1], P-Net [47] and MKD [63]. Results show that on the Retinal-OCT dataset, the proposed SSM outperforms all the state-of-the-art methods, especially the two methods, P-Net and MKD, that utilize a large amount of additional data for the training. In detail, P-Net is a special method for detecting the anomaly in retinal images. It uses structural information of retinal images, which is extracted by a segmentation
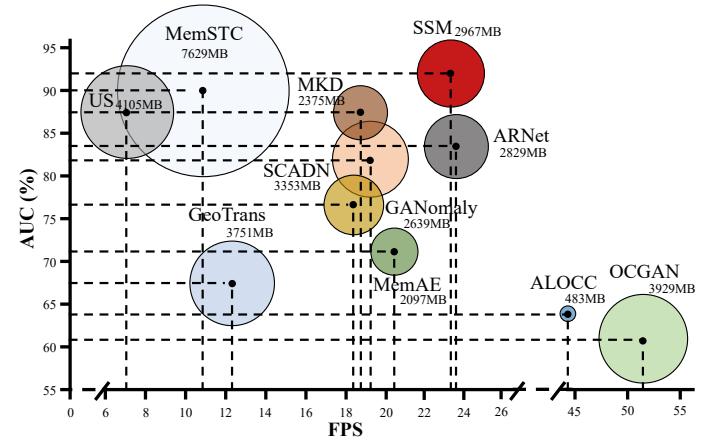


Fig. 6: Comparison of frames per second (FPS) (horizontal coordinates), GPU memory usages (circular sizes) and AUC for anomaly detection (vertical coordinates) of various methods testing on MVTec. The proposed SSM takes up a relatively small GPU memory, and its FPS is relatively higher.

network pre-training with a large amount of additional data and pixel-level annotations [47]. The knowledge distillation-based methods MKD uses ImageNet [82] to pre-train the networks and thus obtains powerful feature representation abilities. The proposed SSM uses only the data in the training set of the Retinal-OCT dataset, without any other additional data. Under this setting, SSM still shows better performance (98.3% AUC) than P-Net (92.9% AUC) and MKD (97.0% AUC), showing the effectiveness of the proposed method.

**Results on MVTec AD dataset.** For the task of anomaly detection on the MVTec AD dataset, Table II shows the corresponding results. We compare the proposed SSM with several state-of-the-arts, including OCGAN [9], ALOCC [8], GeoTrans [17], DAE [42], MemAE [16], GANomaly [50],

TABLE IV: Ablation studies on different losses and architectures of SSM. Results are shown as AUC in % on the Retinal-OCT dataset. The best-performing method is in bold.

| | Losses | | | Architectures | | |
|---|---|---|---|---|---|---|
| $\mathcal{L}_{MSE}$ | $\mathcal{L}_{SSIM}$ | $\mathcal{L}_{GMS}$ | $\mathcal{L}_{MASK}$ | MAM | Refinement | AUC (%) |
| ✓ | | | | | | 94.3 |
| ✓ | | | ✓ | | | 94.7 |
| | ✓ | | ✓ | | | 74.8 |
| | | ✓ | ✓ | | | 94.0 |
| ✓ | ✓ | ✓ | ✓ | | | 96.2 |
| ✓ | ✓ | ✓ | ✓ | ✓ | | 96.8 |
| ✓ | ✓ | ✓ | ✓ | | ✓ | 97.6 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | **98.3** |

TABLE V: Ablation studies for anomaly detection and localization on different grid sizes of masks for SSM. Results are shown as AUC in %. The best-performing method is in bold.

| Grid Sizes of Masks | | | | Detection | | Localization |
|---|---|---|---|---|---|---|
| 4 | 8 | 16 | 32 | Retinal-OCT | MVTec AD | MVTec AD |
| ✓ | ✓ | | | 96.8 | 90.8 | 93.2 |
| | ✓ | ✓ | | 98.0 | 91.4 | 93.4 |
| | | ✓ | ✓ | 94.3 | 91.3 | 93.1 |
| ✓ | ✓ | ✓ | | **98.3** | **92.0** | **93.9** |
| | ✓ | ✓ | ✓ | 94.6 | 91.4 | 93.2 |
| ✓ | ✓ | ✓ | ✓ | 94.9 | 91.9 | 93.8 |

SCADN [22], ARNet [19], US [60] and MKD [63]. As shown in Table II, the proposed SSM achieves the highest mean AUC among all categories (92.0% AUC, 2.0% higher than MemSTC). In 9 out of the 15 categories, SSM outperforms all the other baseline methods. For the other 6 categories, the best performance is achieved by 5 different methods. SSM also achieves the least standard deviation (7.88) for the 15 categories, compared to US (10.92), MKD (7.94), and MemSTC (10.43), which shows that SSM has a good generalizability across different categories.

For the task of anomaly localization, we consider DAE [42], OCGAN [9], MemAE [16], AnoGAN [41], SCADN [22], CNN-dict [78], SMAI [23], GDR [79], MKD [63] and GP [62] as the baselines. The corresponding results are shown in Table III. Overall, results in Table III show that the proposed SSM outperforms other methods on 7 categories and achieves the highest mean AUC among all categories (93.9% AUC, 2.9% higher than GP), showing the effectiveness of anomaly localization of SSM.

**Limitation**. Anomaly regions of different categories vary significantly in attributes such as shape and size. As these attributes are unknown at training, some important hyperparameters, *e.g.*, the mask grid size, cannot be well determined by prior. As a result, we choose to use the same hyper-parameter setting for all categories, which may not guarantee the best performance for every category at the same time. However, results have shown that the proposed SSM is able to yield the best overall performance on all the categories as it achieves the highest mean AUC compared with all the state-of-the-art methods.

**Computational Cost**. We investigate the computational efficiency and the cost of GPU memory for SSM, as well
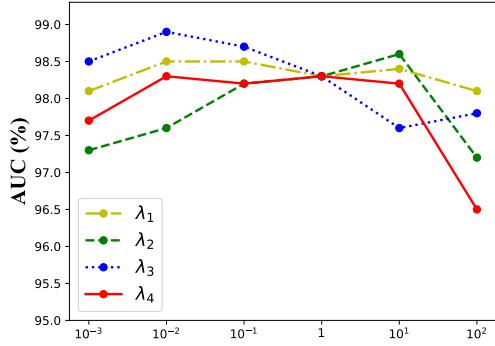


Fig. 7: Sensitivity analysis w.r.t. the hyperparameters for SSM on Retinal-OCT. The AUC for anomaly detection is reported. Best viewed in color.

as several state-of-the-art methods. The results are shown in Figure 6. For all methods, we test 5 times on the MVTec AD dataset with NVIDIA GTX 3090 and record the average FPS and the GPU memory costs. Among all the methods, OCGAN [9] has an advantage in the highest computational efficiency (51.3 fps) but takes up relatively large GPU memory (3929MB). The most state-of-the-art method MemSTC [48] consumes more GPU memory (7629MB) and suffers from a relatively low computational efficiency (10.9 fps). SSM reaches 23.5 fps and takes only 2967MB of GPU memory. Without the mask attention module, SSM reaches a slightly higher speed (24.2fps) and consumes lower GPU memory (2803M). Though during the inference phase, SSM needs to refine the masks with multiple image reconstructions, SSM still reaches a considerable efficiency thanks to its light network structure. Compared with SCADN [22], the most recent state-of-the-art image inpainting-based method, SSM has a much lighter network structure and traverses fewer possible masks during inference. To summarize, with the best performance in AUC, SSM takes up a relatively small GPU memory with high computational efficiency.

*C. Ablation Studies*

We study the contribution of the proposed components of SSM independently.

*1) Loss Functions:* Table IV shows experimental results of ablation studies on Retinal-OCT dataset. We first evaluate the impact of the mask reconstruction loss, $\mathcal{L}_{MASK}$. As $\mathcal{L}_{MASK}$ cannot be use alone to train the model, we compare $\mathcal{L}_{MASK}+\mathcal{L}_{MSE}$ with $\mathcal{L}_{MSE}$ alone, and show that adding the mask reconstruction loss leads to an increase in AUC from 94.3% to 94.7%. We then investigate the influences of the three image reconstruction losses, $\mathcal{L}_{MSE}$, $\mathcal{L}_{SSIM}$ and $\mathcal{L}_{GMS}$. For this set of experiments, $\mathcal{L}_{MASK}$ are always used as default. Among the three losses, $\mathcal{L}_{MSE}$ focuses on the error of every pixel in the image and leads to the highest AUC (94.7%); $\mathcal{L}_{GMS}$, a similarity loss function based on image gradient magnitude, is good at detecting anomalous regions resulted from the roughness and bulge of the object edges (94.0% in AUC); $\mathcal{L}_{SSIM}$, focusing on measuring the structural similarity of images, despite its lowest AUC among the three (74.8%), is good at detecting
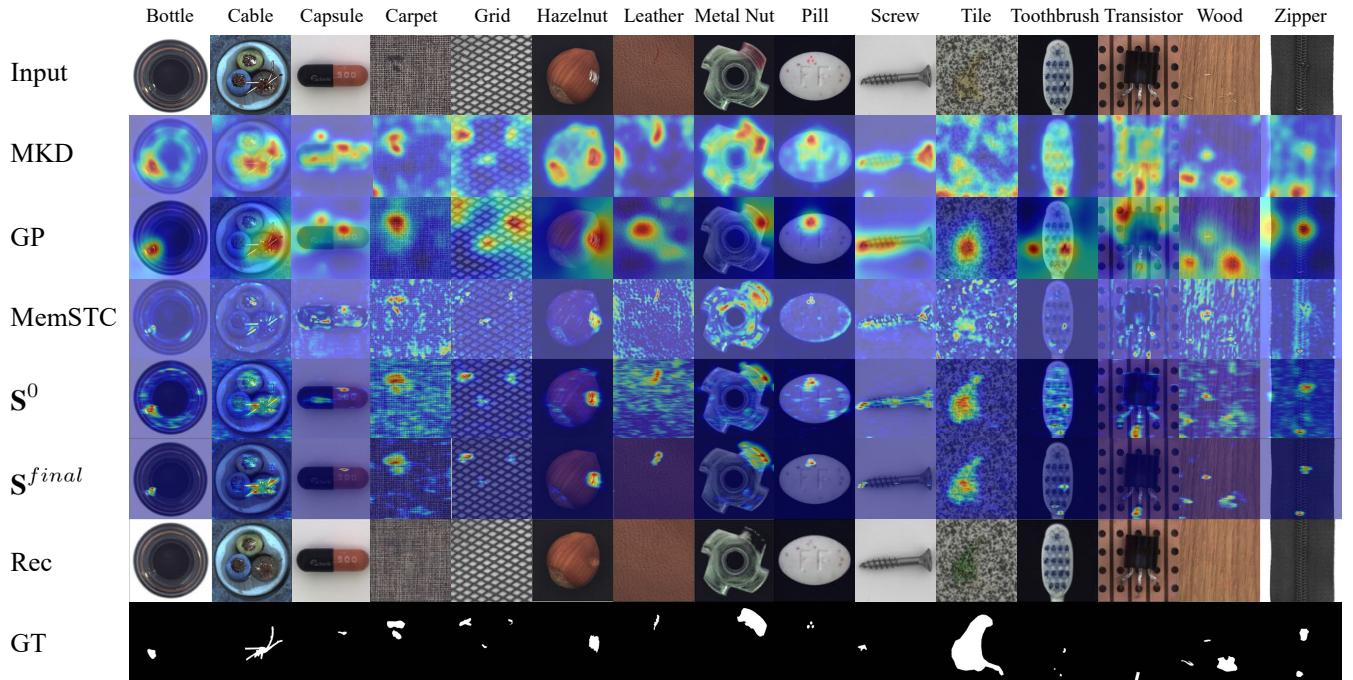
Fig. 8: Qualitative results of anomaly localization of SSM on the MVTec AD dataset for several **difficult cases**, compared with several state-of-the-art methods, including MKD [63], GP [62], and MemSTC [48]. For SSM, $\mathbf{S}^0$ is the anomaly score map at the initialization step. $\mathbf{S}^{final}$ is the final anomaly score map after mask refinement by SSM. Rec is the final reconstructed image by SSM. GT is the ground truth.

structural anomalies. Combining the three image reconstruction losses, together with the mask reconstruction loss, the AUC of anomaly detection can be further improved to 96.2%.

We also note that the baseline method with only the MSE loss has higher AUC (94.3% as shown in Table IV) compared to P-Net [47] (92.9% as shown in Table I), the state-of-the-art reconstruction-based method for anomaly detection, despite the exactly same autoencoder architectures employed. The gain in AUC for the baseline method suggests that the proposed random masking and restoring framework, a form of image inpainting, itself outperforms the whole image reconstruction-based framework which P-Net adopts. A similar result is shown in Table II, where SCADN [22] (a simple image inpainting framework) outperforms the state-of-the-art image reconstruction-based method MemAE [16] for $> 10\%$ AUC on MVTec AD dataset.

*2) Architectures:* The two main important designs of SSM are the mask attention module (MAM) and the progressive mask refinement approach. As one of the self-supervised learning methods, MAM is used to improve semantic feature learning and the representation ability of the latent features. Table IV shows that the mask attention module (MAM) can steadily improve the anomaly detection performance (from 96.2% to 96.8% in AUC). Then, the progressive mask refinement approach (see 'Refinement' in Table IV) can also be used to further improve the anomaly detection performance (from 96.8% to 98.3% in AUC). The higher the value is, the more difficult it is to improve the AUC, which shows that the proposed progressive mask refinement approach is highly effective. Since the anomaly detection approach of

SSM often contains 2-4 iterations, to build a strong baseline, experiments without the progressive mask refinement approach are conducted on reconstruction tasks under 16 randomly generated masks for one test data.

*3) Grid Sizes of Masks:* We discuss the basic grid sizes of the masks used for the model training and the mask initialization during the test. Table V shows the corresponding experimental results on Retinal-OCT and MVTec AD dataset. Results show that with $4 \times 4$, $8 \times 8$ and $16 \times 16$ masks, SSM shows the best anomaly detection and localization performance. For example, we achieve 98.3% AUC for anomaly detection on Retinal-OCT dataset, 92.0% for anomaly detection and 93.9% for anomaly localization on MVTec AD dataset. Experiments containing $32 \times 32$ masks obtain relatively lower AUC, because too large masks greatly increase the difficulty of image reconstruction. We thus use the mask with sizes of $4 \times 4$, $8 \times 8$ and $16 \times 16$ in the following experiments with default. Complete results for each category on MVTec AD are shown in the appendix.

*4) Sensitivity Analysis:* In the previous submission, all experiments are performed under the default setting for SSM, *i.e.*, all individual loss weights $\lambda_i$ ($i = \{1, \cdots, 4\}$) are set to 1. To further study the impact of the weights, we perform a sensitivity analysis for each of the above hyperparameters (fixing the other hyperparameters to 1). As shown in Figure 7, the performance of SSM is not too much sensitive to the choice of the hyperparameters. Although carefully tuning the hyperparameters may be able to lead to a slighter better result, it may also pose the risk of overfitting and also introduce a lot of computational costs. As a result, we leave the hyperparameters to default.
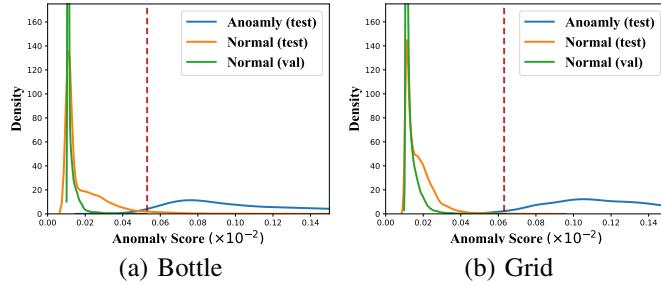
(a) Bottle　　　　　　　　(b) Grid

Fig. 9: The pixel-level anomaly score distributions of (a) bottle and (b) grid on the MVTec AD dataset. The red dotted lines represent the threshold utilized in the mask refinement approach, which is set as the maximum anomaly score in the validation set for each category.
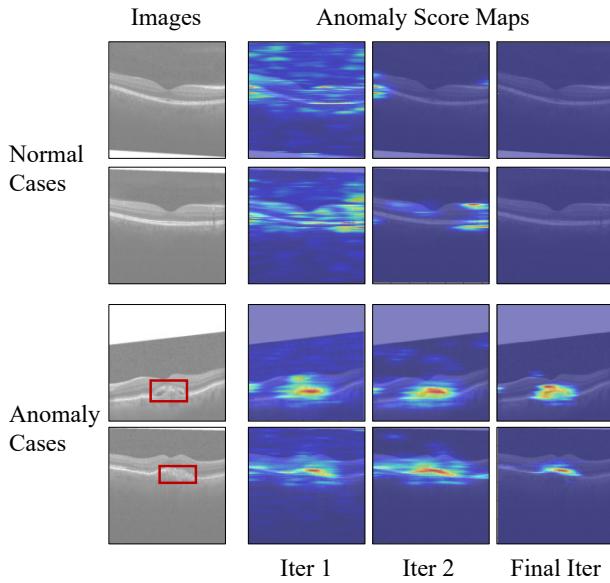


Fig. 10: Qualitative results of anomaly localization with SSM on the Retinal-OCT dataset. During the process of mask refinement, SSM continually narrows the scopes of attention maps to only the anomalous regions.

### D. Visualization Analysis

To analyze how conditional autoencoder and the corresponding progressive mask refinement approach improve the anomaly localization performance, we visualize the results of some hard cases in the MVTec AD dataset and Retinal-OCT dataset to provide qualitative analysis.

**MVTec AD dataset.** Figure 8 shows the visualization of the results from the MVTec AD dataset, including one case for each of the 15 categories, compared with several state-of-the-art methods, including MKD [63], GP [62], and MemSTC [48]. For SSM, the initialized anomaly score map $\mathbf{S}^0$ is wrongly highlighted in many normal regions. Due to some disturbances such as the light intensity, the shadow, or the special placing angle of the object, the model locates the anomaly in part of normal areas after the initialization phase. Fortunately, after the mask refinement, as shown in $\mathbf{S}^{final}$, SSM narrows the scope of attention for anomaly cases in Figure 8 and finally successfully

located the anomalies, which performs significantly better than the other state-of-the-art methods. The visual analysis strongly illustrates the effectiveness of the progressive mask refinement approach.

We also visualize the anomaly score distributions of bottle and grid on MVTec AD dataset in Figure 9. The red dotted lines represent the threshold utilized in the progressive mask refinement approach. The values of the thresholds are set as the maximum anomaly score in the validation set. We can clearly see that the normal and anomalous pixels can be well separated by the thresholds.

**Retinal-OCT dataset.** We show the visualization results for normal and anomaly cases in Retinal-OCT dataset in Figure 10. During the process of mask refinement, SSM continually narrows the scopes of attention maps for both normal and anomaly cases. Especially for the anomaly cases, SSM obtains better predictions in the final iteration, which is matched with the related lesion areas (red bounding in the images). For the normal cases, although SSM focuses on some normal areas wrongly after the initialization step (Iter 1), these areas are eventually be cleared, which helps achieve the correct conclusion (Final Iter). The visual analysis well illustrates the effectiveness of the progressive mask refinement approach.

### V. Conclusion and Future Work

This paper proposes a novel technique named *Self-supervised Masking (SSM)* for unsupervised anomaly detection and localization. A conditional autoencoder is leveraged to learn powerful representations under a larger search space of randomly generated masks in the manner of self-supervised learning. Then we introduce a progressive mask refinement approach to progressively uncover the normal regions and finally locate the anomalous regions. The proposed SSM outperforms state-of-the-arts on multiple anomaly detection benchmarks for both anomaly detection and anomaly localization. Notably, there are still more effective mask refinement approaches to be explored. It is likely to further improve the effectiveness of SSM by better mask refining methods, *e.g.*, improving the updating strategy for multi-scale masks and designing better refinement stopping criterion. The way to locate the targeted regions with a mask refinement module can also be applied to more unsupervised learning fields, opening avenues for future research.

### References

[1] C. Baur, B. Wiestler, S. Albarqouni, and N. Navab, "Deep autoencoding models for unsupervised anomaly segmentation in brain mr images," in *Int. MICCAI Brainlesion Workshop*. Springer, 2018, pp. 161–169.

[2] K. Ouardini, H. Yang, B. Unnikrishnan, M. Romain, C. Garcin, H. Zenati *et al.*, "Towards practical unsupervised anomaly detection on retinal images," in *Domain Adaptation and Representation Transfer and Medical Image Learning with Less Labels and Imperfect Data*. Springer, 2019, pp. 225–234.

[3] N. Shvetsova, B. Bakker, I. Fedulova, H. Schulz, and D. V. Dylov, "Anomaly detection in medical imaging with deep perceptual autoencoders," *IEEE Access*, vol. 9, pp. 118 571–118 583, 2021.

[4] J. Zhang, Y. Xie, Z. Liao, G. Pang, J. Verjans, W. Li, Z. Sun, J. He, and C. S. Yi Li, "Viral pneumonia screening on chest x-ray images using confidence-aware anomaly detection," *IEEE Trans. Med. Imag.*, vol. 40, no. 3, pp. 879–890, 2021.

[5] T. Matsubara, R. Tachibana, and K. Uehara, "Anomaly mach. component detection by deep generative model with unregularized score," in *Int. Joint Conf. Neural Netw.* IEEE, 2018, pp. 1–8.

[6] C. Phua, V. Lee, K. Smith, and R. Gayler, "A comprehensive survey of data mining-based fraud detection research," in *Int. Conf. on Intelligent Computation Technology and Automation*, vol. 1, 2010, pp. 50–53.

[7] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song, "Robust physical-world attacks on deep learning visual classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2018, pp. 1625–1634.

[8] M. Sabokrou, M. Khalooei, M. Fathy, and E. Adeli, "Adversarially learned one-class classifier for novelty detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2018, pp. 3379–3388.

[9] P. Perera, R. Nallapati, and B. Xiang, "Ocgan: One-class novelty detection using gans with constrained latent representations," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2019, pp. 2898–2906.

[10] Z. Zhang, S. Chen, and L. Sun, "P-kdgan: Progressive knowledge distillation with gans for one-class novelty detection," in *Int. Joint Conf. Artif. Intelligence*, 2020, pp. 3237–3243.

[11] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, "Estimating the support of a high-dimensional distribution," *Neural Computation*, vol. 13, no. 7, pp. 1443–1471, 2001.

[12] K. Lee, H. Lee, K. Lee, and J. Shin, "Training confidence-calibrated classifiers for detecting out-of-distribution samples," in *Int. Conf. Learn. Representations*, 2018.

[13] L. Ruff, R. Vandermeulen, N. Goernitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, and M. Kloft, "Deep one-class classification," in *Proc. Int. Conf. Mach. Learn.* PMLR, 2018, pp. 4393–4402.

[14] B. Zong, Q. Song, M. R. Min, W. Cheng, C. Lumezanu *et al.*, "Deep autoencoding gaussian mixture model for unsupervised anomaly detection," in *Int. Conf. Learn. Representations*, 2018.

[15] L. Deecke, R. Vandermeulen, L. Ruff, S. Mandt, and M. Kloft, "Image anomaly detection with generative adversarial networks," in *Proc. Joint Eur. Conf. Mach. Learn. Knowl. Discov. Databases.* Springer, 2018, pp. 3–17.

[16] D. Gong, L. Liu, V. Le, B. Saha, M. R. Mansour, S. Venkatesh, and A. v. d. Hengel, "Memorizing normality to detect anomaly: Memory-augmented deep autoencoder for unsupervised anomaly detection," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019, pp. 1705–1714.

[17] I. Golan and R. El-Yaniv, "Deep anomaly detection using geometric transformations," in *Proc. Neural Inf. Process. Syst.*, vol. 31, 2018, pp. 9758–9769.

[18] S. Wang, Y. Zeng, X. Liu, E. Zhu, J. Yin, C. Xu, and M. Kloft, "Effective end-to-end unsupervised outlier detection via inlier priority of discriminative network," in *Proc. Neural Inf. Process. Syst.*, 2019, pp. 5960–5973.

[19] F. Ye, C. Huang, J. Cao, M. Li, Y. Zhang, and C. Lu, "Attribute restoration framework for anomaly detection," *IEEE Trans. Multimedia*, vol. 24, pp. 116–127, 2022.

[20] D. S. Kermany, M. Goldbaum, W. Cai, C. C. Valentim, H. Liang, S. L. Baxter, A. McKeown, G. Yang, X. Wu, F. Yan *et al.*, "Identifying medical diagnoses and treatable diseases by image-based deep learning," *Cell*, vol. 172, no. 5, pp. 1122–1131, 2018.

[21] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger, "Mvtec ad–a comprehensive real-world dataset for unsupervised anomaly detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2019, pp. 9592–9600.

[22] X. Yan, H. Zhang, X. Xu, X. Hu, and P.-A. Heng, "Learning semantic context from normal samples for unsupervised anomaly detection," in *Proc. AAAI Conf. Artif. Intelligence*, vol. 35, no. 4, 2021, pp. 3110–3118.

[23] Z. Li, N. Li, K. Jiang, Z. Ma, X. Wei, X. Hong, and Y. Gong, "Superpixel masking and inpainting for self-supervised anomaly detection," in *British Mach. Vis. Conf.*, 2020, pp. 7–10.

[24] B. R. Kiran, D. M. Thomas, and R. Parakkal, "An overview of deep learning based methods for unsupervised and semi-supervised anomaly detection in videos," *J. Imag.*, vol. 4, no. 2, p. 36, 2018.

[25] W. Chu, H. Xue, C. Yao, and D. Cai, "Sparse coding guided spatiotemporal feature learning for abnormal event detection in large videos," *IEEE Trans. Multimedia*, vol. 21, no. 1, pp. 246–255, 2018.

[26] K. Xu, X. Jiang, and T. Sun, "Anomaly detection based on stacked sparse coding with intraframe classification strategy," *IEEE Trans. Multimedia*, vol. 20, no. 5, pp. 1062–1074, 2018.

[27] K. Xu, T. Sun, and X. Jiang, "Video anomaly detection and localization based on an adaptive intra-frame classification network," *IEEE Trans. Multimedia*, vol. 22, no. 2, pp. 394–406, 2019.

[28] M. Sabokrou, M. Fathy, and M. Hoseini, "Video anomaly detection and localisation based on the sparsity and reconstruction error of autoencoder," *Electronics Letters*, vol. 52, no. 13, pp. 1122–1124, 2016.

[29] M. Sabokrou, M. Pourreza, M. Fayyaz, R. Entezari, M. Fathy, J. Gall, and E. Adeli, "Avid: Adversarial visual irregularity detection," in *Asian Conf. Comput. Vis.* Springer, 2018, pp. 488–505.

[30] M. Sabokrou, M. Fathy, G. Zhao, and E. Adeli, "Deep end-to-end one-class classifier," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 32, no. 2, pp. 675–684, 2020.

[31] M. Sabokrou, M. Fayyaz, M. Fathy, Z. Moayed, and R. Klette, "Deep-anomaly: Fully convolutional neural network for fast anomaly detection in crowded scenes," *Comput. Vis. Image Understanding*, vol. 172, pp. 88–97, 2018.

[32] M. Sabokrou, M. Fayyaz, M. Fathy, and R. Klette, "Deep-cascade: Cascading 3d deep neural networks for fast anomaly detection and localization in crowded scenes," *IEEE Trans. Image Process.*, vol. 26, no. 4, pp. 1992–2004, 2017.

[33] A. Bansal, S. Karan, S. Gaurav, C. Rama, and D. Ajay, "Zero-shot object detection," in *Proc. Eur. Conf. Comput. Vis.* Springer, 2018, pp. 384–400.

[34] C. Yan, X. Chang, M. Luo, H. Liu, X. Zhang, and Q. Zheng, "Semantics-guided contrastive network for zero-shot object detection," *IEEE Trans. Pattern Anal. Mach. Intelligence*, 2022.

[35] E. Eskin, "Anomaly detection over noisy data using learned probability distributions," in *Proc. Int. Conf. Mach. Learn.*, 2000, pp. 255–262.

[36] K. Yamanishi, J. I. Takeuchi, G. Williams, and P. Milne, "On-line unsupervised outlier detection using finite mixtures with discounting learning algorithms," *Data Mining & Knowl. Discovery*, vol. 8, no. 3, pp. 275–300, 2004.

[37] M. Rahmani and G. K. Atia, "Coherence pursuit: Fast, simple, and robust principal component analysis," *IEEE Trans. Signal Process.*, vol. 65, no. 23, pp. 6260–6275, 2017.

[38] H. Xu, C. Caramanis, and S. Sanghavi, "Robust pca via outlier pursuit," *IEEE Trans. Inf. Theory*, vol. 58, no. 5, pp. 3047–3064, 2012.

[39] J. An and S. Cho, "Variational autoencoder based anomaly detection using reconstruction probability," *Special Lecture IE*, vol. 2, no. 1, pp. 1–18, 2015.

[40] Y. Xia, X. Cao, F. Wen, G. Hua, and J. Sun, "Learning discriminative reconstructions for unsupervised outlier removal," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2015, pp. 1511–1519.

[41] T. Schlegl, P. Seeböck, S. M. Waldstein, U. Schmidt-Erfurth, and G. Langs, "Unsupervised anomaly detection with generative adversarial networks to guide marker discovery," in *Int. Conf. Inf. Process. Med. Imag.* Springer, 2017, pp. 146–157.

[42] M. Sakurada and T. Yairi, "Anomaly detection using autoencoders with nonlinear dimensionality reduction," in *Mlsda Workshop Mach. Learn. Sensory Data Anal.*, 2014, pp. 4–11.

[43] M. Nicolau, J. McDermott, and V. L. Cao, "A hybrid autoencoder and density estimation model for anomaly detection," in *Int. Conf. Parallel Prob. Solving from Nat.* Springer, 2016, pp. 717–726.

[44] M. Nicolau, J. McDermott *et al.*, "One-class classification for anomaly detection with kernel density estimation and genetic programming," in *Eur. Conf. Genetic Programming.* Springer, 2016, pp. 3–18.

[45] C. Zhou and R. C. Paffenroth, "Anomaly detection with robust deep autoencoders," in *Proc. ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, 2017, pp. 665–674.

[46] V. Zavrtanik, M. Kristan, and D. Skočaj, "Reconstruction by inpainting for visual anomaly detection," *Pattern Recog.*, vol. 112, p. 107706, 2021.

[47] K. Zhou, Y. Xiao, J. Yang, J. Cheng, W. Liu, W. Luo *et al.*, "Encoding structure-texture relation with p-net for anomaly detection in retinal images," in *Proc. Eur. Conf. Comput. Vis*, 2020, pp. 360–377.

[48] K. Zhou, J. Li, Y. Xiao, J. Yang, J. Cheng, W. Liu, W. Luo, J. Liu, and S. Gao, "Memorizing structure-texture correspondence for image anomaly detection," *IEEE Trans. Neural Netw. Learn. Syst.*, pp. 1–15, 2021.

[49] H. Zenati, C. S. Foo, B. Lecouat, G. Manek, and V. R. Chandrasekhar, "Efficient gan-based anomaly detection," in *Int. Conf. Learn. Representations Workshop*, 2018.

[50] S. Akçay, A. Atapour-Abarghouei, and T. P. Breckon, "Skip-ganomaly: Skip connected and adversarially trained encoder-decoder anomaly detection," in *Int. Joint Conf. Neural Netw.*, 2019, pp. 1–8.

[51] X. Wang, Y. Du, S. Lin, P. Cui, Y. Shen, and Y. Yang, "advae: A self-adversarial variational autoencoder with gaussian anomaly prior knowledge for anomaly detection," *Knowledge-Based Syst.*, vol. 190, p. 105187, 2020.

[52] M. Noroozi and P. Favaro, "Unsupervised learning of visual representations by solving jigsaw puzzles," in *Proc. Eur. Conf. Comput. Vis.* Springer, 2016, pp. 69–84.

[53] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," in *Int. Conf. Learn. Representations*, 2019.

[54] C. Doersch, A. Gupta, and A. A. Efros, "Unsupervised visual representation learning by context prediction," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2015, pp. 1422–1430.

[55] D. Yuan, X. Chang, P.-Y. Huang, Q. Liu, and Z. He, "Self-supervised deep correlation tracking," *IEEE Trans. Image Process.*, vol. 30, pp. 976–985, 2021.

[56] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2016, pp. 2536–2544.

[57] M. Everingham, M. A. S. Eslami, V. L. Gool, K. I. C. Williams, M. J. Winn, and A. Zisserman, "The pascal visual object classes challenge: A retrospective," *Int. J. Comput. Vis*, vol. 111, no. 1, pp. 98–136, 2015.

[58] E. Denton, S. Gross, and R. Fergus, "Semi-supervised learning with context-conditional generative adversarial networks," *arXiv preprint arXiv:1611.06430*, 2016.

[59] S. Jenni and P. Favaro, "Self-supervised feature learning by learning to spot artifacts," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2018, pp. 2733–2742.

[60] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger, "Uninformed students: Student-teacher anomaly detection with discriminative latent embeddings," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2020, pp. 4183–4192.

[61] C.-L. Li, K. Sohn, J. Yoon, and T. Pfister, "Cutpaste: Self-supervised learning for anomaly detection and localization," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2021, pp. 9664–9674.

[62] S. Wang, L. Wu, L. Cui, and Y. Shen, "Glancing at the patch: Anomaly localization with global and local feature comparison," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2021, pp. 254–263.

[63] M. Salehi, N. Sadjadi, S. Baselizadeh, M. H. Rohban, and H. R. Rabiee, "Multiresolution knowledge distillation for anomaly detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2021, pp. 14 902–14 912.

[64] Q. Yu, J. Zhang, H. Zhang, Y. Wang, Z. Lin, N. Xu, Y. Bai, and A. Yuille, "Mask guided matting via progressive refinement network," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2021, pp. 1154–1163.

[65] Q. Zhang, J. Lai, Z. Feng, and X. Xie, "Seeing like a human: Asynchronous learning with dynamic progressive refinement for person re-identification," *IEEE Trans. Image Process.*, vol. 31, pp. 352–365, 2021.

[66] Q. Liu and Z. Wang, "Progressive boundary refinement network for temporal action detection," in *Proc. AAAI Conf. Artif. Intelligence*, vol. 34, no. 7, 2020, pp. 11 612–11 619.

[67] X. Huang and Y. Peng, "Tpckt: Two-level progressive cross-media knowledge transfer," *IEEE Trans. Multimedia*, vol. 21, no. 11, pp. 2850–2862, 2019.

[68] M. H. Fayek, L. Cavedon, and R. H. Wu, "Progressive learning: A deep learning framework for continual learning," *Neural Netw.*, vol. 128, pp. 345–357, 2020.

[69] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang, "Free-form image inpainting with gated convolution," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019, pp. 4471–4480.

[70] W. Xue, L. Zhang, X. Mou, and A. C. Bovik, "Gradient magnitude similarity deviation: A highly efficient perceptual image quality index," *IEEE Trans. Image Process.*, vol. 23, no. 2, pp. 684–695, 2013.

[71] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Trans. Image Process.*, vol. 13, no. 4, pp. 600–612, 2004.

[72] X. Zhu, L. Zhang, L. Zhang, X. Liu, Y. Shen, and S. Zhao, "Gan-based image super-resolution with a novel quality loss," *Mathematical Problems in Engineering*, pp. 1–12, 2020.

[73] A. Shahsavari, S. Ranjbari, and T. Khatibi, "Proposing a novel cascade ensemble super resolution generative adversarial network (cesr-gan) method for the reconstruction of super-resolution skin lesion images," *Informatics in Medicine Unlocked*, p. 100628, 2021.

[74] T. Schlegl, P. Seeböck, S. M. Waldstein, G. Langs, and U. Schmidt-Erfurth, "f-anogan: Fast unsupervised anomaly detection with generative adversarial networks," *Med. image Anal.*, vol. 54, pp. 30–44, 2019.

[75] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2017, pp. 1125–1134.

[76] X. Han, X. Chen, and L.-P. Liu, "Gan ensemble for anomaly detection," in *Proc. AAAI Conf. Artif. Intelligence*, vol. 35, no. 5, 2021, pp. 4090–4097.

[77] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2017, pp. 2223–2232.

[78] P. Napoletano, F. Piccoli, and R. Schettini, "Anomaly detection in nanofibrous materials by cnn-based self-similarity," *Sensors*, vol. 18, no. 1, p. 209, 2018.

[79] D. Dehaene, O. Frigo, S. Combrexelle, and P. Eline, "Iterative energy-based projection on a normal data manifold for anomaly localization," in *Int. Conf. Learn. Representations*, 2020.

[80] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Proc. Int. Conf. Med. Image Comput. Comput.-Assist. Interv.* Springer, 2015, pp. 234–241.

[81] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Int. Conf. Learn. Representations*, 2015.

[82] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *Int. J. Comput. Vis*, vol. 115, no. 3, pp. 211–252, 2015.

**Chaoqin Huang** (Student Member, IEEE) received his B.E. degree in computer science from Shanghai Jiao Tong University (SJTU), Shanghai, China, in 2019. He has been working towards a Ph.D. degree at the Cooperative Meidianet Innovation Center, Shanghai Jiao Tong University under the supervision of Prof. Ya Zhang, since 2019. He is also an intern at Shanghai AI Laboratory. His research interests include anomaly detection, computer vision, and machine learning.

**Qinwei Xu** (Student Member, IEEE) received the B.E. degree in photoelectric information science and engineering from the University of Electronic Science and Technology of China, in 2018. He is currently pursuing the Ph.D. degree with the Cooperative Medianet Innovation Center, Shanghai Jiao Tong University, Shanghai, China, under the supervision of Prof. Ya Zhang. He is also an intern at Shanghai AI Laboratory. His research interests include deep learning and computer vision.

**Yanfeng Wang** (Member, IEEE) received the B.S. degree from PLA Information Engineering University, and the master's and Ph.D. degrees in business management from Shanghai Jiao Tong University. He is currently the Vice Director of Cooperative Medianet Innovation Center and also the Vice Dean of the School of Electrical and Information Engineering, Shanghai Jiao Tong University. His research interest is mainly on media big data, emerging commercial applications of information technology, and technology transfer.

**Yu Wang** (Member, IEEE) received the B.E. degree from Huazhong University of Science and Technology, Wuhan, China, in 2009, and the M.Sc. degree in Communications and Signal processing and the Ph.D. degree in Signal Processing, both from Imperial College, London, U.K. in 2010 and 2015, respectively. He was a Senior Research Associate with Machine Intelligence Laboratory, Cambridge University Engineering Department, when he was a key member of the ALTA and IARPA MATERIAL projects. Since December 2020, he has been a Associate Professor with Cooperative Medianet Innovation Center, Shanghai Jiao Tong University. His current research interests include machine learning, audio signal processing, speech recognition, spoken language processing and multi-modal signal processing. He is a Member of IEEE and ISCA.

**Ya Zhang** (Member, IEEE) is currently a professor at the Cooperative Medianet Innovation Center, Shanghai Jiao Tong University. Her research interest is mainly in machine learning with applications to multimedia and healthcare. Dr. Zhang has published more than 100 refereed papers in prestigious international conferences and journals. She has won several best paper awards of international journals and conferences, and directed one Outstanding Doctorate Dissertations awarded by Chinese Association for Artificial Intelligence.

APPENDIX

TABLE VI: Complete results of ablation studies on different grid sizes of masks for SSM about **anomaly detection** on the MVTec AD dataset. Results are listed as AUC in % and are marked individually for each class. An average score over all classes is also reported in the last row. The best-performing method is in bold.

| Category | Grid Sizes of Masks | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | [4] | [8] | [16] | [32] | [4, 8] | [8, 16] | [16, 32] | [4, 8, 16] | [8, 16, 32] | [4, 8, 16, 32] |
| Bottle | 98.2 | 98.7 | 99.5 | 99.5 | 99.7 | 99.8 | 100 | 99.9 | 99.9 | 100 |
| Cable | 59.1 | 60.8 | 67.3 | 76.6 | 74.3 | 74.8 | 77.5 | 77.3 | 76.3 | 77.2 |
| Capsule | 87.2 | 86.5 | 90.0 | 90.5 | 89.7 | 91 | 91.6 | 91.4 | 92.5 | 91.7 |
| Carpet | 67.7 | 67.9 | 76.8 | 72.7 | 72.9 | 74.8 | 73.9 | 76.3 | 73.3 | 73.4 |
| Grid | 94.2 | 100 | 100 | 99.9 | 100 | 100 | 100 | 100 | 100 | 100 |
| Hazelnut | 81.6 | 88.5 | 92.2 | 95.4 | 88.1 | 91.5 | 93.5 | 91.5 | 95.9 | 93.5 |
| Leather | 99.0 | 100 | 99.8 | 99.4 | 100 | 99.7 | 99.8 | 99.9 | 99.3 | 99.0 |
| Metal Nut | 90.2 | 74.2 | 90.3 | 88.8 | 76.7 | 86.9 | 89.6 | 88.7 | 91.7 | 90.2 |
| Pill | 94.4 | 92.1 | 89.7 | 83.9 | 94.2 | 89.6 | 82.8 | 89.1 | 82.1 | 85.6 |
| Screw | 80.5 | 86.5 | 83.8 | 82.1 | 87.5 | 85.3 | 81.0 | 85.0 | 82.3 | 83.5 |
| Tile | 97.7 | 96.1 | 94.9 | 89.6 | 95.9 | 93.3 | 91.0 | 94.4 | 89.6 | 96.6 |
| Toothbrush | 96.7 | 98.3 | 100 | 100 | 99.4 | 100 | 100 | 100 | 100 | 100 |
| Transistor | 74.0 | 82.5 | 86.0 | 93.5 | 84.9 | 90.3 | 95.5 | 91.0 | 93.5 | 92.7 |
| Wood | 98.0 | 98.2 | 93.8 | 90.6 | 98.7 | 94.5 | 94.5 | 95.9 | 94.9 | 95.2 |
| Zipper | 99.9 | 99.9 | 99.9 | 99.6 | 99.9 | 100 | 99.6 | 99.9 | 99.6 | 99.8 |
| Mean | 87.5 | 89.3 | 90.9 | 90.8 | 90.8 | 91.4 | 91.3 | **92.0** | 91.4 | 91.9 |

TABLE VII: Complete results of ablation studies on different grid sizes of masks for SSM about **anomaly localization** on the MVTec AD dataset. Results are listed as AUC in % and are marked individually for each class. An average score over all classes is also reported in the last row. The best-performing method is in bold.

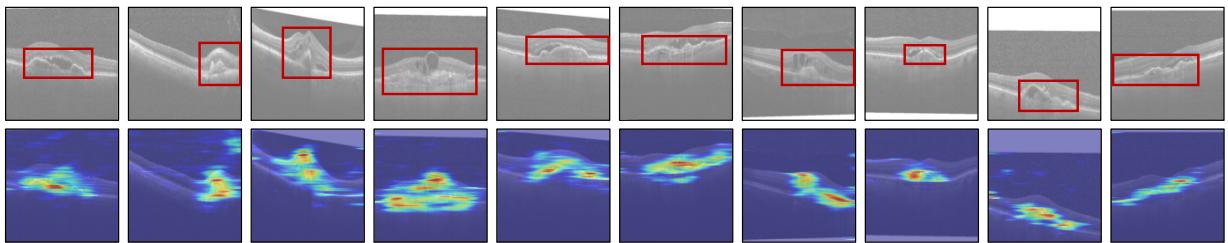| Category | Grid Sizes of Masks | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | [4] | [8] | [16] | [32] | [4, 8] | [8, 16] | [16, 32] | [4, 8, 16] | [8, 16, 32] | [4, 8, 16, 32] |
| Bottle | 92.2 | 93.8 | 95.6 | 96.1 | 94.1 | 96.0 | 96.4 | 95.9 | 96.5 | 96.6 |
| Cable | 81.1 | 78.9 | 75.8 | 81.2 | 78.9 | 79.2 | 81.3 | 82.1 | 77.4 | 80.9 |
| Capsule | 93.4 | 95.1 | 98.0 | 98.4 | 97.3 | 98.1 | 98.7 | 98.4 | 98.5 | 98.5 |
| Carpet | 90.4 | 87.1 | 94.2 | 92.0 | 92.8 | 91.5 | 92.1 | 94.4 | 92.9 | 92.9 |
| Grid | 98.5 | 99.0 | 98.9 | 98.7 | 99.0 | 99.0 | 98.8 | 99.0 | 98.8 | 98.8 |
| Hazelnut | 96.3 | 96.9 | 97.7 | 98.2 | 96.9 | 97.8 | 98.1 | 97.4 | 98.0 | 98.1 |
| Leather | 99.5 | 99.7 | 99.6 | 99.6 | 99.6 | 99.6 | 99.6 | 99.6 | 99.6 | 99.6 |
| Metal Nut | 77.8 | 81.4 | 89.1 | 93.2 | 83.7 | 89.5 | 92.4 | 89.6 | 93.2 | 92.8 |
| Pill | 99.1 | 98.8 | 98.0 | 96.5 | 98.8 | 97.6 | 96.7 | 97.8 | 96.6 | 96.6 |
| Screw | 98.0 | 98.9 | 98.9 | 98.8 | 99.0 | 99.0 | 98.9 | 98.9 | 99.0 | 99.0 |
| Tile | 96.6 | 93.5 | 88.4 | 81.9 | 93.4 | 88.7 | 79.1 | 90.2 | 83.8 | 87.7 |
| Toothbrush | 98.5 | 98.8 | 98.9 | 98.9 | 98.8 | 98.8 | 98.9 | 98.9 | 98.9 | 98.9 |
| Transistor | 76.4 | 77.7 | 79.8 | 83.2 | 78.5 | 80.1 | 83.7 | 80.1 | 82.9 | 83.1 |
| Wood | 89.0 | 88.7 | 87.0 | 84.7 | 88.4 | 87.1 | 83.7 | 86.9 | 83.7 | 84.1 |
| Zipper | 99.0 | 99.0 | 99.0 | 98.9 | 99.0 | 99.0 | 98.9 | 99.0 | 98.9 | 98.9 |
| Mean | 92.8 | 92.5 | 93.3 | 93.4 | 93.2 | 93.4 | 93.1 | **93.9** | 93.2 | 93.8 |



Fig. 11: More qualitative results of anomaly localization with SSM on the Retinal-OCT dataset.

To better show the main contributions in this paper, we summarize the differences between the traditional image reconstruction-based method and the image inpainting-based method in Table VIII.

| | Without Masking | With Masking |
|---|---|---|
| Framework | Image Reconstruction | Image Restoration / Image Inpainting |
| Inputs | Entire Image | Image + Mask |
| Assumption | The model is enforced to learn regularities of normal data to minimize reconstruction errors; anomalies are difficult to be reconstructed and thus have large reconstruction errors. | The model is enforced to predict the missing information from its surrounding normal regions; the difference between the original masked region and its corresponding restoration result is significant for anomalous regions. |
| Challenge | Compressing the original image and then reconstructing it. | The dual tasks of *masking* possible anomalies and *restoring* the masked regions. |
| Pros/Cons | <ul><li>No need for mask design and mask selection.</li><li>Learned representation may focus on low level details.</li></ul> | <ul><li>Need careful design and selection of the mask.</li><li>If the mask lies in the anomalous regions correctly, the image inpainting-based method performs significantly better than the image reconstruction-based method.</li></ul> |

TABLE VIII: Summarization of the key contribution of the proposed method.