# Visual Anomaly Detection Via Partition Memory Bank Module and Error Estimation

Peng Xing, Zechao Li

*Abstract*—Reconstruction method based on the memory module for visual anomaly detection attempts to narrow the reconstruction error for normal samples while enlarging it for anomalous samples. Unfortunately, the existing memory module is not fully applicable to the anomaly detection task, and the reconstruction error of the anomaly samples remains small. Towards this end, this work proposes a new unsupervised visual anomaly detection method to jointly learn effective normal features and eliminate unfavorable reconstruction errors. Specifically, a novel Partition Memory Bank (PMB) module is proposed to effectively learn and store detailed features with semantic integrity of normal samples. It develops a new partition mechanism and a unique query generation method to preserve the context information and then improves the learning ability of the memory module. The proposed PMB and the skip connection are alternatively explored to make the reconstruction of abnormal samples worse. To obtain more precise anomaly localization results and solve the problem of cumulative reconstruction error, a novel Histogram Error Estimation module is proposed to adaptively eliminate the unfavorable errors by the histogram of the difference image. It improves the anomaly localization performance without increasing the cost. To evaluate the effectiveness of the proposed method for anomaly detection and localization, extensive experiments are conducted on three widely-used anomaly detection datasets. The encouraging performance of the proposed method compared to the recent approaches based on the memory module demonstrates its superiority.

*Index Terms*—Anomaly detection, Partition memory bank, Histogram error estimation module.

## I. INTRODUCTION

VISUAL anomaly detection aims to detect abnormal data that are different from normal visual data [1]–[3], which has shown great potential in a variety of applications, such as industrial anomaly detection [4]–[8] and medical diagnosis [9]. Since the abnormal patterns are diverse and the occurrence frequency of visual anomaly data is much lower than the normal ones, it is impractical to obtain sufficient abnormal training samples with different abnormal patterns. Consequently, it is challenging to successfully detect the anomaly data by using only normal training data.

Some methods have been studied to address the anomaly detection task. Recently, in [12]–[14], Autoencoder (AE) [15] is introduced for unsupervised anomaly detection to detect abnormal samples based on reconstruction errors [12], [16]–[18]. These methods expect AE to reconstruct normal samples with only small reconstruction errors and reconstruct abnormal samples with large reconstruction errors. It is well

P. Xing, Z. li are with the School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing 21094, China. E-mail: xingp_ng@njust.edu.cn, zechao.li@njust.edu.cn. (Corresponding Author: Zechao Li)
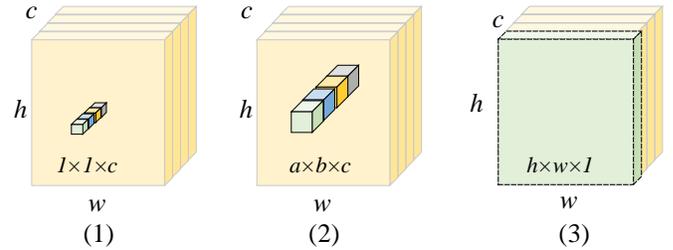


Fig. 1. Illustration of the query generation methods for different memory modules. Figures (1) [10] and (2) [11] represent the logical query generated from the same position in the feature map, with different channels. Figure (3) represents the proposed query. It preserves the critical features of original feature map learned by the convolutional network.

known that AE has a strong reconstruction capability [10], [19] in many computer vision tasks. [20]–[23]. However, the powerful reconstruction ability is not beneficial for anomaly detection. Even for unlearned anomalous samples, AE can still reconstruct them, resulting in small reconstruction error for unlearned anomalous samples [10]. Therefore, the reconstruction error of AE reconstruction alone is not sufficient to solve the anomaly detection challenge. Some works propose more complex self-supervised tasks, e.g., image-colorization [24]. The complex self-supervised task makes AE learn normal semantic information to recover the original image. For the unlearned abnormal samples, the abnormal semantics cannot be recovered, making the reconstruction error larger. However, with tiny differences in semantic information between abnormal and normal samples, these methods can not detect the abnormal samples quite well.

To make the reconstruction errors of abnormal samples larger than the ones of normal samples, some methods [10], [11], [19], [25] combine AE and the traditional memory modules. The output of the encoder is used as the query of a memory module, and the read of this memory module is used as the input of the decoder. Then, the memory module stores the normal features from the encoder output. When the abnormal feature queries the memory module, the output is normal feature. Therefore, the reconstruction error of abnormal sample becomes large. However, due to the query generation method, the existing memory modules do not actually store normal features, but rather logical information. As shown in Figure 1 (1) and Figure 1 (2), the query generated by the traditional method is a logical reorganization of features at the same location but from different channels of the original feature map [19] [11]. A single channel feature is extracted by the same convolutional kernel by the sliding window
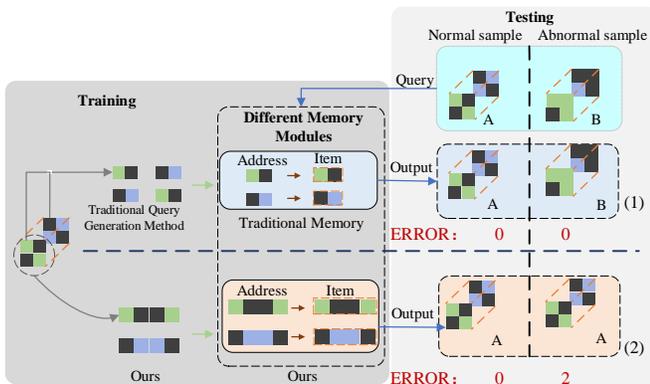
Fig. 2. The differences in query generation method and memory item between the proposed method and the traditional method. "ERROR" represents reconstruction error. The memory modules are both learned from the reconstruction task of category A to category A using reconstruction loss update. However, when the anomalous sample B is different from the normal sample A, the traditional methods can still successfully reconstruct B ("ERROR" is 0), while the proposed method fails to reconstruct unlearned category samples ("ERROR" is 2), which makes it successfully detect the anomalous sample B.

method from the feature map of the previous layer, which contains similar critical feature information. Nevertheless, the queries generated by the logical reorganization of the feature map destroy the critical feature learned by the convolutional network, allowing the memory module to learn the reorganized logical information. When anomalous features as query address memory module, the traditional module will fit the anomalous feature reads in logical relationship to support image reconstruction. Therefore, it makes the detection and localization of anomalous samples fail with small reconstruction errors. Figure 2 (1) shows an example of detection failure for anomalous sample B, which learns a memory module by the traditional query generation method (Figure 1 (1)).

When facing tiny region anomaly samples, there is another aspect to consider - the effect of cumulative reconstruction error. Introducing the memory module makes it impossible for AE to achieve zero-error reconstruction even for normal regions. Due to the cumulative error, the abnormal scores of some normal samples can surpass the abnormal sample scores of small abnormal regions (e.g., small scratches), leading to incorrect detection and localization results.

Moreover, Existing methods [10], [19] adopt a scheme, which concatenates the outputs of last encoder and the reads of the memory module as the input of the decoder. However, it can not well explore the storage from the memory module since encoder branch can provide sufficient features for reconstruction. Furthermore, these methods place the memory module after the last layer of the encoder focusing more on high-level features, which are not applicable for anomaly localization requiring low-level features. Consequently, all anomaly detection methods based on AE unable to reach the pixel-level anomaly localization.

Toward this end, this paper proposes a novel AE joint memory module network, as shown in Figure 3. It utilizes a Partition Memory Bank (PMB) module to improve the abilities of learning and storing normal features, as well as a

Histogram Error Estimation module to adaptively eliminate the effects of cumulative reconstruction errors caused by memory module. To preserve the critical feature information learned from the original feature map, the PMB module develops a novel query generation method that utilizes each channel feature to generate a query, as shown in Figure 1 (3). It allows the memory module to learn normal semantic features rather than reorganization logical features. When anomalous features as query address the PMB module, it reads the stored normal features in a semantic relationship to get normal features. Therefore, the reconstruction error becomes larger. Figure 2(2) shows an example of a successful detection where the PMB module causes an abnormal sample B to be reconstructed as a normal sample A (the reconstruction error is large enough to be successfully detected as abnormal). Meanwhile, PMB needs to make the normal region reconstruction error as small as possible. A partition mechanism is proposed to further improve the expressiveness of the memory module, which uses multiple memory units to independently store the normal features of different partitions. Hence, PMB module enables the small reconstruction error for normal regions while the large error for the abnormal regions, which benefits the detection of the abnormal samples. The Histogram Error Estimation module utilizes the reconstruction error maps to construct histograms. Then, it is used to adaptively estimate the reconstruction error for the normal region of each image to obtain a more effective corrected error map. More importantly, it is a non-parametric method and does not consume additional resources.

To make the PMB module put more attention on low-level features for more precise localization, a new scheme is developed. We directly utilize skip connections as input of decoder on high-level features and memory module reads as input on low-level features. The reads of PMB and the outputs of the skip connections are separately utilized at different layers rather than concatenated at the same layer, which allows the memory module to learn more detailed features of normal samples. Experimental results on widely-used datasets demonstrate that the proposed method obtains competitive results.

The main contributions are summarized as follows:

(1) This paper proposes a novel Partition Memory Bank (PMB) module. The unique partition mechanism and a query generation method of PMB can effectively learn and store normal features, and achieve excellent results in anomaly detection task.

(2) To address the challenge of more accurate anomaly localization, a new non-parametric Histogram Error Estimation module is developed to eliminate the cumulative reconstruction error, which can obtain better anomaly detection results and anomaly localization maps.

(3) The AE, PMB module, and Histogram Error Estimation module are jointly explored for their optimal compatibility. Experiments are conducted on three benchmark datasets, which shows that the proposed method can effectively solve the problem of successful reconstruction of abnormal samples.
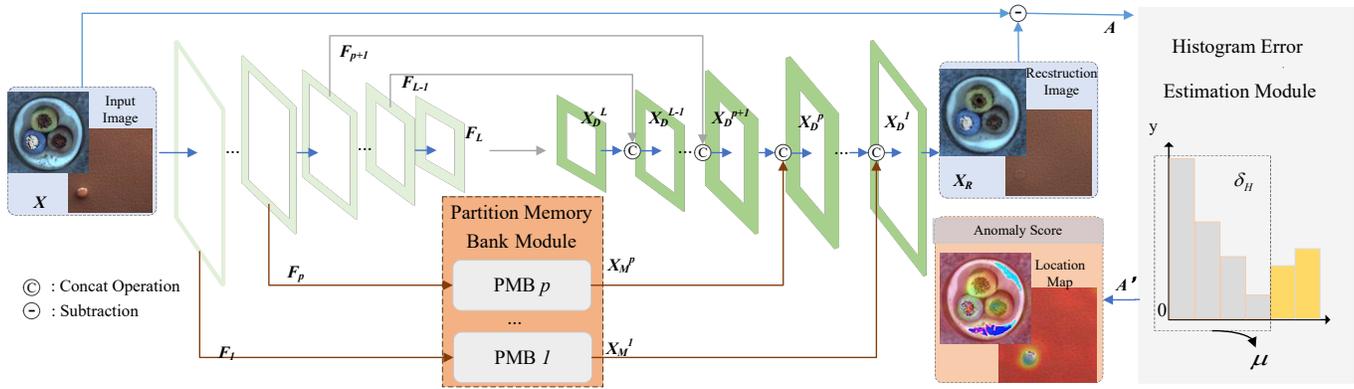
Fig. 3. Illustration of the proposed framework. It contains the AE, the PMB module, and the histogram error estimation module. First, the $L$-layer feature map is extracted by the encoder module in the AE structure. The PMB module stores the first $P$-layer low-level features and subsequently uses skip connections. The reconstructed image $X_R$ forms an error map $A$ with the original image $X$. Second, the Histogram Error Estimation module uses $A$ to build a histogram, estimate and eliminate $\mu$, and obtain a corrected error map $A'$. $A'$ is used for anomaly detection and localization.

## II. RELATED WORK

### A. Unsupervised Visual Anomaly Detection

For unsupervised anomaly detection, the detection model is trained by only normal samples, which is used to detect abnormal samples and localize abnormal regions [2], [13]. Traditional classifiers are introduced for anomaly detection, such as SVM [26] and OC-SVM [27]. With the popularity of deep learning, researchers proposed deep one-class methods such as DSVDD [28] and OCNN [29]. The common idea of these methods is to learn the decision boundary of normal samples. For example, DSVDD learns a spherical discriminant plane to improve discrimination efficiency. There are also other methods that model the distribution features of normal samples, such as MRF [30], MDT [31] and GMM [2], [17]. However, these methods are less effective to deal with high dimensional data.

Because of the outstanding performance of transfer learning in other computer vision tasks, the pretrained network is introduced to the field of anomaly detection [6], [32], [33]. In [6], since the student network only learns latent representations of normal samples, the degree of sample abnormality is measured by the difference between latent feature representations from multiple student and teacher networks. However, they rely on additional training datasets and huge resource consumption. Some researches have proposed the employment of GAN [34] to solve the anomaly detection problem. The generator in GAN is used to learn normal sample features and discriminator to identify anomalies [35]–[37]. However, GAN may be unstable and tends to derive unsatisfactory results in the actual situation.

Recently, reconstruction-based methods have become popular in anomaly detection, which expects to poorly reconstruct abnormal samples and enlarges the gap of the reconstruction error of abnormal samples and normal samples [12]–[14], [24], [38]–[41]. AE [15] is used to reconstruct samples. For example, AE-SSIM [14] attempts to reconstruct normal samples directly using AE, while ARFAD [24] learns the reconstruction features of normal samples through self-supervised tasks such as rotation. However, these methods successfully reconstruct

not only normal samples but also abnormal samples, which leads to limited performance.

Different from the above methods, this work proposes a new PMB module to reconstruct successfully normal samples and poorly abnormal samples.

### B. Memory Network

The memory network was initially adapted to the field of natural language processing [20], [21], [42]–[44]. Recently, Some existing works apply memory networks to anomaly detection [19], [25]. Gong et al. [10] proposes MemAE, which uses a memory module to reduce the reconstruction ability of AE for anomalous samples. However, due to its simple memory module and the shortcomings of the combination method, it is difficult to be effective for tiny anomaly sample detection and localization tasks. Park et al. [19] proposes a combination network of AE and memory module for video anomaly detection. It constructs an update method for the memory items of the memory module and concatenates the outputs of AE and memory module as the input of the decoder. However, it is not suitable for the localization challenges of complex scenes because it focuses more on high-level features. Second, the output of AE as the input of the decoder results in the limitation of features stored in the memory module. Wang et al. [25] introduces an additional evaluation network on the MemAE framework as a discriminator to detect whether a sample is anomalous or not. However, it is difficult to achieve anomalous region localization with discriminators.

The most related study to the proposed method is DAAD. [11]. It uses the method in Figure 1 (2) to generate memory query of size $a \times b \times c$ ($a \leq h$, $b \leq w$) and uses a discriminator to identify whether the samples are anomalous or not. DAAD uses larger query (generally $8 \times 8 \times c$) to make anomalies and normals that do not share the same block pattern. However, it destroys the integrity of the features learned by the convolutional network, resulting in the memory module learning reorganized logical information. Therefore, abnormal features can also be logically fitted to anomalies such that reconstruction errors are small. The proposed PMB
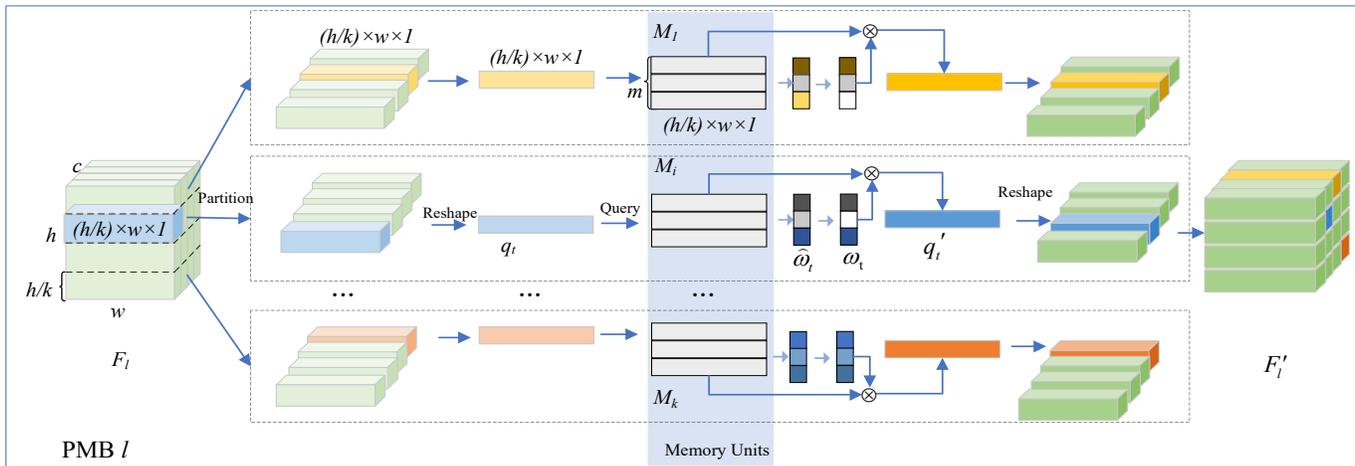
Fig. 4. The framework of the proposed PMB module. For each layer of feature map $F_l$, the partition mechanism causes the feature map to be divided into separate partitions. Each partition has a separate memory unit to store normal features. Each query addresses the memory unit in an attentional way to obtain the valid normal feature reads as input to the decoder.

module employs novel query of size $h \times w \times 1$ (Figure 1(3)) which enables the memory module to store normal feature learned from the convolutional network. Therefore, all queries can address the PMB module to get normal features. Unlike DAAD, it is not degrade the reconstruction capability of AE, but make the memory module unable to address abnormal features. In addition, DAAD introduces discriminator with adversarial loss to solve the problem of incorrect discrimination caused by cumulative errors, which results in additional resource consumption and the failure of anomaly localization. In this paper, we not only ensure that the normal features are well reconstructed, but also the anomaly localization results can be derived from the error map without additional resources.

## III. THE PROPOSED METHOD

### A. Overview

This work proposes to solve the unsupervised anomaly detection problem based on the reconstruction scheme. The overview of the proposed framework is shown in Figure 3. The proposed framework contains three basic structures: 1) an AE, 2) a PMB module, and 3) a Histogram Error Estimation module.

Given the input sample $X$, the encoder first extracts the multi-scale latent representations $F = \{F_1, F_2, \ldots, F_L\}$. To make full utilization of the representation information provided by the memory modules and allow the memory module to focus on low-level features, we utilize the output of the PMBs at layer 1 to layer $p$ of the encoder and the output of the subsequent $p+1$ to $L$ layers with the skip connection. Therefore, $p$ PMBs are designed to learn and store feature maps of different scales independently. The reads of the proposed PMB module are represented by $X_M = \{X_M^1, X_M^2, \ldots, X_M^p\}$. Then, the channel concatenation between $X_M$ and output of previous decoder $X_D$ as the input of the next decoder. The

reconstructed image $X_R$ is obtained after the decoder. The original difference image $A$ is obtained by using Eq. 1,

$$A = |\mathrm{X} - X_R|. \tag{1}$$

Besides, the Histogram Error Estimation module is developed to eliminate reconstruction errors in normal regions to obtain a more accurate corrected error map $A'$ for anomaly localization and detection. The anomaly scores and anomaly localization maps of the samples are obtained based on $A'$.

### B. Partition Memory Bank Module

As is shown in Figure 4, The PMB module is proposed to learn and store the potential feature representations of normal samples. It achieves the challenge of anomalous sample reconstruction error becoming large by exploring the proposed joint partition mechanism and a new query generation method.

To generate the query with critical feature information, a novel query generation method is developed, as shown in Figure 1 (3). It generates queries directly from each channel of the feature map. Thus, each query contains semantic information for a single channel, allowing the PMB module to store only normal semantic information. In addition, to enhance the learning ability and expressiveness of the memory module and to make the reconstruction error of normal regions small, a partition mechanism is proposed. It introduces multiple local units to store the normal features of different regions separately. Each memory unit is capable of storing detailed feature information for individual regions. Multiple memory units make the expressiveness of the memory module enhanced.

Specifically, for each feature map $F_l$ with the size of $h \times w \times c$ from $F$, it is partitioned by using the proposed query generation approach shown in Figure 4. Along the '$h$' dimension of the feature map, $F_l$ is partitioned into $k$ regions with the size of $\frac{h}{k} \times w \times c$, each of which is stored and queried by a separate memory unit. The feature map of individual channel in each partition is flattened to

one $\frac{h}{k} \times w$ dimensional vector. Thus, we can get query sets $Q = \{q_t | t \in [1, c]\}$ in each partition. Each PMB contains $k$ memory units $M = \{M_1, M_2, \ldots, M_k\}$, each memory unit contains $m$ memory items, and the size of the memory items is the same as the query $q_t$. Finally, as shown in Fig. 4 (4), suppose $q_t$ is a query for $i$-th partition of $F_l$, $q_t$ and the corresponding memory items in $M_i$ are first normalized to improve the accuracy of the attention weight. The normalization operations are as follows:

$$\widehat{q_t} = norm\left(q_t\right), \widehat{M_i^j} = norm\left(M_i^j\right), j \in [1, m], \quad (2)$$

where $norm(\cdot)$ represents the $l_2$-norm. $\widehat{q_t}$ represents the normalized query $q_t$. $\widehat{M_i}$ represents the memory items after normalization. Then the similarity between $\widehat{q_t}$ and the memory item $\widehat{M_i}$ is used to calculate the attention weight $\widehat{\omega_t}$ as follows:

$$\widehat{\omega_t^j} = \frac{exp(<\widehat{q_t}, \widehat{M_i^j}>)}{\sum\limits_{u=1}^{m} exp(<\widehat{q_t}, \widehat{M_u^j}>)} \quad (3)$$

Here $< \cdot, \cdot >$ denotes the cosine similarity. To avoid too small weight of the query, this paper proposes to filter the weights by introducing a threshold, which can alleviate the reconstruction with most of memory items [10]. Consequently, when abnormal features are used as queries, the reads are difficult to fit and the abnormal regions are then difficult to be reconstructed.

$$\omega_t^j = max\left(0, \widehat{\omega_t^j} - \delta_m\right), \quad (4)$$

where $\delta_m$ represents the threshold. By using the memory module, we can obtain a new feature vector $q_t'$ for each $q_t$ by using the output of the memory unit $M_i$.

$$q_t' = \omega_t \cdot M_i \quad (5)$$

Then a new feature map $F'$ can be obtained by re-arranging all $q_t'$ and fed into the decoder to generate the reconstructed image $X_R$. It is worth noting that the stability of the reconstructed network can be improved by using $q_t'$ since the value of $q_t'$ does not change much.

### C. Histogram Error Estimation Module

For each image, the original error map $A$ is obtained by using Eq. 1, in which the value of each pixel represents the reconstruction error. In the proposed approach, the reconstruction error comes from the proposed PMB module and AE. The PMB module learns and stores the normal features, which enables to guarantee that the reconstruction errors of normal regions are obviously smaller than ones of abnormal regions. However, due to the presence of the memory module, AE is difficult to achieve zero-error reconstruction in the normal region. These accumulative reconstruction errors can result in false detection from the original error map. Therefore, we need to further mitigate the impact of accumulative reconstruction errors on detection and localization performance.

In this paper, we explore a simple, yet effective Histogram Error Estimation module to estimate the error $\mu$. As shown in Figure 3, the histogram of error map is first constructed. Due to the small reconstruction error of normal region, the pixels in normal region always lie on the left side of histogram. Therefore, we can quickly find normal pixels in the histogram. Then, a percentage $\delta_H$ of small reconstruction errors are chosen to estimate $\mu$ as shown in the dashed box of the histogram in Figure 3. The average value of errors of these selected pixels is utilized to estimate $\mu$. Then, the corrected difference image $A'$ is obtained as follows:

$$A' = \begin{cases} A - \mu; & A > \mu \\ 0; & A \leq \mu \end{cases} \quad (6)$$

The anomaly score and localization map of the samples are derived from $A'$. $A'$ is adopted as the localization map, while the anomaly score is calculated by the sum of the error values of the corrected difference image $A'$.

### D. Loss Function

To train the proposed model, the quality of the reconstructed image and the intuitive feeling of the human visual system are jointly explored. The reconstruction loss $\mathcal{L}_{MSE}$ and structural similarity index measure (SSIM) [55] loss $\mathcal{L}_{SSIM}$ are introducesd , which are defined as:

$$\mathcal{L}_{MSE} = \|X - X_R\|_2 \quad (7)$$

$$\mathcal{L}_{SSIM} = \frac{1}{h \times w} \sum_{i=1}^{h} \sum_{j=1}^{w} (1 - SSIM_s^{(i,j)}(X, X_R)) \quad (8)$$

Here $h$ and $w$ represent the height and width of the image. $SSIM_s^{(i,j)}(X, X_R)$ represents structural similarity index between the image $X$ and image $X_R$ at the center $(i, j)$ position with kernel size of s.

This paper proposes to jointly explore the above loss to obtain the total training loss as follows:

$$\mathcal{L} = \mathcal{L}_{MSE} + \lambda_{SSIM} \cdot \mathcal{L}_{SSIM}, \quad (9)$$

where $\lambda_{SSIM}$ is the hyperparameter of the model that measures the importance of the SSIM loss.

## IV. EXPERIMENTS

This section conducts extensive experiments to verify the effectiveness of the proposed method with the comparison to several state-of-the-art methods.

### A. Experimental Setting

**Datasets.** To validate the effectiveness of the proposed method, experiments are conducted on three widely-used datasets.

(1) **MVTec AD** [4]. The MVTec AD dataset contains 5354 high-resolution color images in 15 different categories. Among them, There are 5 categories for textured images, such as "wood" or "leather". The other 10 categories contain non-textured objects, such as "cable". The challenge of this dataset comes from the fact that normal images and abnormal images come from the same category. The difference between the

TABLE I
COMPARISON OF THE PROPOSED METHOD AND ANOMALY DETECTION METHODS ON MVTec AD DATASET WITH AUROC (%). THE BEST RESULTS ARE MARKED IN BOLD. '✓' MEANS THAT THE METHOD ADDITIONALLY INTRODUCES LARGE DATASETS SUCH AS IMAGENET.

| | Method | Extra Datasets | Bottle | Hazelnut | Capsule | Metal Nut | Leather | Pill | Wood | Carpet | Tile | Grid | Cable | Transistor | Toothbrush | Screw | Zipper | Mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pre-trained | CNN-Dict [45] | ✓ | 78 | 72 | 84 | 82 | 87 | 68 | 91 | 72 | 93 | 59 | 79 | 66 | 77 | 87 | 76 | 78 |
| | MKDAD [32] | ✓ | 99.4 | 98.4 | 80.5 | 73.6 | 95.1 | 82.7 | 94.3 | 79.3 | 91.6 | 78.0 | 89.2 | 85.6 | 92.2 | 83.3 | 93.2 | 87.7 |
| | SPADE [46] | ✓ | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 85.5 |
| | FAVAE [47] | ✓ | 99.9 | 99.3 | 80.4 | 85.2 | 67.5 | 82.1 | 94.8 | 67.1 | 80.5 | 97.0 | 95.0 | 93.2 | 95.8 | 83.7 | 97.2 | 87.9 |
| | Cutpaste [48] | ✓ | 100 | 99.7 | 94.3 | 98.7 | 100 | 91.3 | 99.8 | 100 | 98.9 | 98.8 | 93.9 | 95.6 | 92.8 | 86.0 | 99.9 | **96.6** |
| AE | CAVGA-$D_w$ [49] | ✓ | 93 | 90 | 89 | 81 | 80 | 93 | 89 | 80 | 81 | 79 | 86 | 80 | 96 | 79 | 95 | 86 |
| | Puzzle-AE [50] | - | 94:2 | 91:2 | 66:9 | 66:3 | 72:9 | 71:6 | 89:5 | 65:7 | 65:5 | 75:3 | 87:9 | 85:9 | 97:8 | 57:8 | 75:7 | 77.6 |
| | ARFAD [24] | - | 94.1 | 85.5 | 68.1 | 66.7 | 86.2 | 78.6 | 92.3 | 70.6 | 73.5 | 88.3 | 83.2 | 84.3 | 100.0 | 100.0 | 87.6 | 83.9 |
| | MemAE [10] | - | 95.4 | 89.1 | 83.1 | 53.7 | 61.1 | 88.3 | 95.4 | 45.4 | 63.0 | 94.6 | 69.4 | 79.3 | 97.2 | 99.2 | 87.1 | 80.2 |
| | AESc [51] | - | 98.0 | 94.0 | 74.0 | 73 | 89.0 | 84.0 | 95.0 | 89.0 | 99.0 | 97.0 | 89.0 | 91.0 | 100 | 74.0 | 94.0 | 89.0 |
| | **OURS** | - | 95.2 | 99.4 | 82.3 | 84.5 | 94.5 | 86.1 | 100.0 | 93.1 | 97.2 | 97.1 | 85.6 | 92.4 | 95.8 | 97.0 | 77.3 | **91.8** |

TABLE II
AUROC OF THE PROPOSED METHOD AND THE COMPARED METHODS FOR ANOMALY DETECTION ON MNIST.

| Method | DSVDD [28] | CapsNetPP [52] | OCGAN [53] | LSA [54] | MemAE [10] | OCSVM [27] | ARAE [18] | **OURS** | MKDAD [32] |
|---|---|---|---|---|---|---|---|---|---|
| Pre-trained model | - | - | - | - | - | - | - | - | ✓ |
| AUROC | 94.8 | 97.7 | 97.5 | 97.5 | 97.5 | 96.0 | 97.5 | 98.1 | **99.35** |

abnormal sample and the normal sample is subtle, such as scratches on the "leather". Each category contains a dataset including only normal images for training and a dataset including normal and various abnormal images for testing.

(2) **MNIST** [56]. The dataset contains 10 different types of low-resolution grayscale images. A total of 70,000 images. 60,000 pictures for training, and 10,000 pictures for testing. In anomaly detection, a category is used as a normal sample and the remaining 9 categories are used as anomaly samples (called out-of-distribution detection [57]).

(3) **Retinal-OCT** [58] It is a medical dataset, which contains normal (healthy) retinal CT images and 3 categories of abnormal (damaged) retinal CT images.

**Implementation Details.** In experiments, each image is resized into the size of $128 \times 128$. The UNet [59]-like network is utilized as the autoencoder in this paper. The learning rate is set to 0.002. For MVTec AD, $L$ is set to 5 and $p$ is set to 4. The number of $k$ in each PMB is set to $(32, 16, 16, 8)$ for five texture categories and $(16, 16, 16, 8)$ for other object categories. Accordingly, the number of memory items $m$ is $(150, 200, 200, 200)$. For MNIST, $L$ and $p$ are both set to 3. The number of partitions $k$ is set to $(1, 1, 1)$ and $m$ is set to $(10, 10, 10)$, respectively. For Retinal-OCT, we set $L$ and $p$ to 4 and 3, respectively. $k$ in each PMB is set to $(2, 2, 2)$, respectively. The number of memory items $m$ is set to $(10, 10, 10)$, respectively. The batch size is set to 20 for the MVTec AD and Retinal-OCT datasets, as well as 64 for MNIST. For the hyper-parameters in the proposed formulation, follow [25], $\delta_m$ is set to $\frac{1}{m}$, $\delta_H$ and $\lambda_{SSIM}$ are set to 50% and 0.1, respectively.

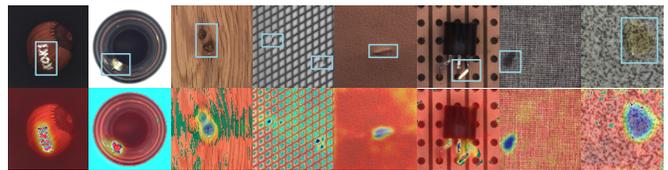**Evaluation.** Following [19], [36], the area under the curve



Fig. 5. Illustration of the localization heatmap. Anomaly images are shown in the top row and the corresponding weighted heatmaps are presented in the bottom row. The locations of abnormal regions are denoted by the rectangular box.

(AUC) of the receiver operating characteristic (ROC) of the image level and pixel level are used to measure the performance of model.

### B. Experimental Results

**Anomaly detection results.** The anomaly detection results on these three datasets and the comparison with the state-of-the-art methods are shown in Table I, II and III, respectively.

(1) **MVTec AD.** For the MVTec AD dataset, this paper compares the state-of-the-art methods with distribution-based methods and AE methods. The distribution-based approach used pre-trained models with ImageNet dataset [64], such as CNN-Dict [45], MKDAD [32], SPADE [46], FAVAE [47], Cutpaste [48]. Among them, MKDAD uses the VGG-16 [65] model pre-trained on ImageNet, SPADE uses the ResNet-18 [66] model pre-trained on ImageNet, and CutPaste uses the EfficientNet (B4) [67] model pre-trained on ImageNet. AE methods include CAVGA-$D_w$ [49], Puzzle-AE [50], ARFAD [24], memAE [10], AESc [51]. For example, Puzzle-AE introduces the complex self-supervised task of puzzle reduction,

TABLE III
AUROC OF THE PROPOSED METHOD AND THE COMPARED METHODS FOR ANOMALY DETECTION ON RETINAL-OCT.

| Method | DSVDD [28] | AnoGan [35] | VAE-GAN [60] | Cycle-GAN [61] | GANomaly [36] | P-Net [62] | MKDAD [32] | **OURS** |
|---|---|---|---|---|---|---|---|---|
| pre-trained model | - | - | - | - | - | - | ✓ | - |
| AUROC | 74.4 | 84.81 | 90.64 | 87.39 | 91.96 | 92.88 | 97.01 | **98.00** |

TABLE IV
AUROC OF THE PROPOSED METHOD AND THE COMPARED METHODS FOR ANOMALY LOCALIZATION ON MVTEC AD.

| Method | Extra Datasets | Bottle | Hazelnut | Capsule | MetalNut | Leather | Pill | Wood | Carpet | Tile | Grid | Cable | Transistor | Toothbrush | Screw | Zipper | Mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| AE-SSIM [14] | - | 93 | 97 | 94 | 89 | 78 | 91 | 73 | 87 | 59 | 94 | 82 | 90 | 92 | 96 | 88 | 87 |
| AE-$L_2$ [14] | - | 86 | 95 | 88 | 86 | 75 | 85 | 73 | 59 | 51 | 90 | 86 | 86 | 93 | 96 | 77 | 82 |
| AnoGAN [35] | - | 86 | 87 | 84 | 76 | 64 | 87 | 62 | 54 | 50 | 58 | 78 | 80 | 90 | 80 | 78 | 74 |
| CNN-Dict [45] | ✓ | 78 | 72 | 84 | 82 | 87 | 68 | 91 | 72 | 93 | 59 | 79 | 66 | 77 | 87 | 76 | 78 |
| MKDAD [32] | ✓ | 96.3 | 94.6 | 95.8 | 86.4 | 98.1 | 89.6 | 84.8 | 95.6 | 82.8 | 91.8 | 82.4 | 76.5 | 96.1 | 95.9 | 93.9 | 90.7 |
| CAVGA-$D_u$ [49] | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 85.0 |
| UTAD [63] | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 90.0 |
| CAVGA-$D_w$ [49] | ✓ | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 92.0 |
| **OURS** | - | 93.7 | 92.5 | 92.1 | 85.8 | 96.7 | 91.1 | 86.5 | 92.3 | 90.7 | 94.3 | 94.2 | 80.8 | 97.5 | 97.7 | 95.4 | **92.1** |

ARFAD introduces the self-supervised task of rotation and coloration, and MemAE introduces the traditional memory module.

From Table I, the proposed method achieves mean AUROC 91.8%. It is lower than Cutpaste due to the fact that Cutpaste uses a better pre-trained model and uses forged anomalous samples for training. However, the proposed method outperforms the remaining methods using pre-trained models by 4%-10%, validating the effectiveness of the AE combined with the PMB module. Among all AE methods, our method achieves state-of-the-art performance. It outperforms the CAVAG-$D_w$ method that introduces real anomalous samples for training by 6%, validating that our method also performs very well without relying on anomalous samples. It is 13% and 8% higher than Puzzle-AE and ARFAD methods, indicating that our method can achieve excellent performance without designing complex image preprocessing operations. It outperforms the traditional memory model MemAE method by 11%, validating that the proposed PMB module stores normal features more efficiently.

(2) **MNIST.** For the MNIST dataset, this paper compares the state-of-the-art methods including DSVDD [28], CapsNetPP [52], OCGAN [53], LSA [54], MemAE [10], OCSVM [27], ARAE [18], , and MKDAD [32]. From Table II, it can be found that the proposed method achieves mean AUROC 98.1% [57] in the out-of-distribution detection setting. The advanced performance shows the superiority of the proposed method. It is slightly lower than the pre-trained method MKDAD method but far outperforms the remain distribution methods. Besides, compared with the traditional memory module MemAE, the proposed method outperforms MemAE, which validates the effectiveness of the novel query method and partition mechanism.

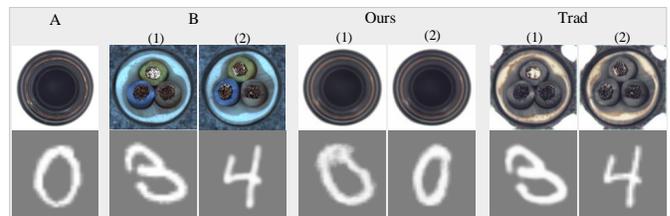(3) **Retinal-OCT.** In the medical anomaly detection task, the



Fig. 6. Illustrations of reconstructed images by using the proposed PMB module or the traditional memory module. 'A' represents normal samples and 'B' represents abnormal samples, 'Ours' shows the test results of the PMB module and 'Trad' shows the test results of the traditional memory module. The reconstruction error of the proposed PMB module ('Ours') reconstructing anomalous sample 'B' is very large and can be effectively detected. However, the traditional module has a small reconstruction error for abnormal samples('Trad').

proposed method achieves significant performance advantages. We compare with DSVDD [28], AnoGan [35], VAE-GAN [60], Cycle-GAN [61], GANomaly [36], P-Net [62], and MKDAD [32] on Retinal-OCT dataset. The results are shown in Table III. It can be found that the proposed method achieves the state-of-the-art performance with AUROC of 98.0%. It not only outperforms recent methods using GAN networks [60], [61], but even outperforms the MKDAD method. Our method does not rely on additional training data inductive bias specific to the pre-trained model. It is more suitable for medical image anomaly detection.

**Anomaly localization results.** As far as we know, none of the recent AE (including memory module) anomaly detection methods achieve pixel-level anomaly detection (anomaly localization task), which may be due to the fact that these methods weak AE reconstruction capability resulting in cannot guarantee small reconstruction error for normal region. However, the proposed method does not weaken the reconstruction
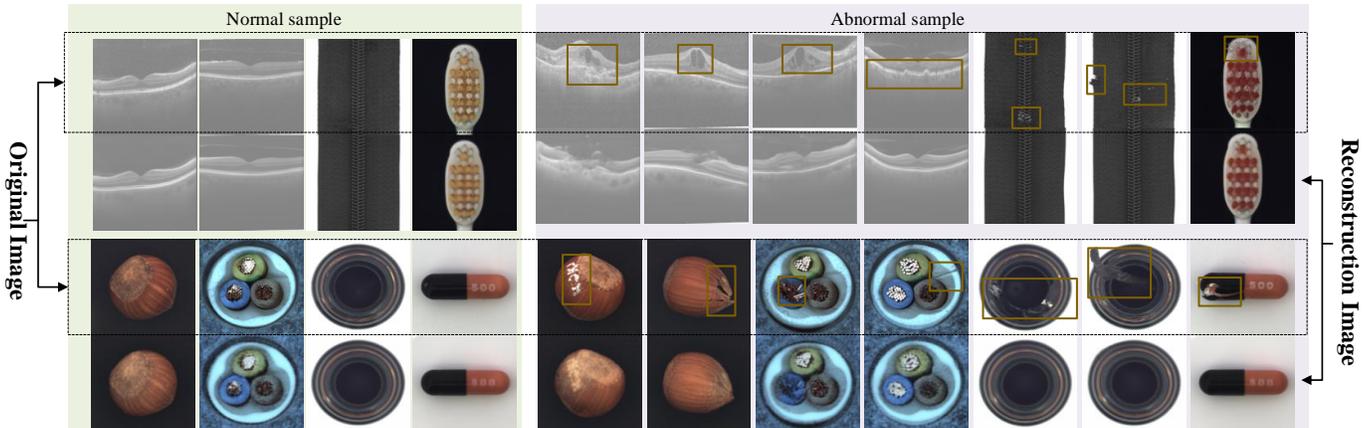
Fig. 7. Some examples of reconstruction of abnormal and normal samples. The original images and the reconstructed images are listed in the top row and the bottom row, respectively.

capability of the model but makes the anomalous features unaddressable from the PMB module thereby increasing the reconstruction error of the anomalous samples. Moreover, the unfavorable error is eliminated by using the Histogram Error Estimation module to further enlarge the difference. The results of anomaly localization on the MVTec AD dataset are shown in Table IV. As a result, our method achieves 92.1% performance in anomaly localization.

Compared with AE-based methods such as AE-SSIM [14] and AE-$L_2$ [14], the proposed method gains 5%-10% improvement. Because the proposed PMB module can store the learned normal sample features, effectively making the anomalous region reconstruction error larger. Compared with the two-stage method UTAD [63], the proposed method is simpler and more effective in training, and better localization results are obtained by exploring the corrected differential image $A'$. Compared with the weakly supervised method CAVGA-$D_w$ [49] which uses additional abnormal region annotation information, the proposed method achieves better results. It shows that the proposed method outperforms the method with weakly supervised settings in the unsupervised setting, fully validating its superiority. To further demonstrate the localization performance of the proposed method, some examples of anomaly localization are illustrated in Figure 5. It can be observed that the proposed method can well solve the problem of abnormal region localization.

### C. Ablation study

**The impact of the different memory models.** To validate the advantages of the proposed PMB module over the traditional memory module for anomaly detection, we used the OOD detection setting to demonstrate the reconstruction content: one category as the normal samples and the remaining categories as abnormal samples. We conducted visualization studies in two datasets (MVTec AD and MNIST) separately. The results are illustrated in Figure 6. For example, by using category '0' as the normal samples for training, the samples '3' and '4' were reconstructed during the test. It can be observed that since the PMB module cannot successfully reconstruct
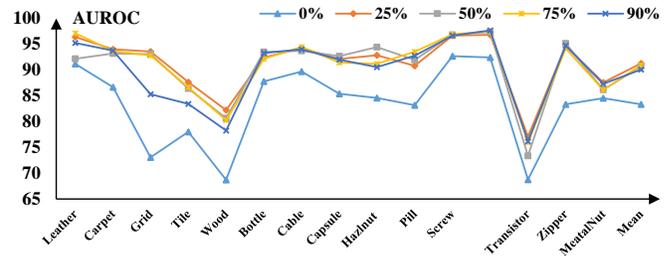


Fig. 8. The influence of different values of $\delta_H$ for anomaly localization.

the samples of class '3' and '4', the reconstruction results are shown as class '0'. The abnormal samples produce a large reconstruction error. In contrast, the traditional memory module can successfully reconstruct the class '3' and '4' samples with a small reconstruction error. Therefore, it is difficult to solve the anomaly detection problem by using the traditional memory module to learn the logical pixels, while the proposed PMB module can ensure the storage of normal feature information and use it for normal feature reconstruction. It can be found the similar results in the MVTec AD dataset, the PMB module has a large reconstruction error for abnormal samples.

As shown in Figure 7, we show the reconstruction results of the PMB module for more abnormal and normal samples. In the 'Normal sample' on the left side of Figure 7, the reconstruction error is small. In the 'Abnormal Sample' on the right, the reconstruction error of the anomalous region shown in the rectangular box is large, while the rest of the errors are small. It indicates that the PMB module does not weaken the reconstruction capability of AE, but makes it impossible for anomalous features to successfully address the abnormal feature output, while normal features can successfully address the normal feature output.

**Compare with different AE methods and memory module methods.** To show the advantages of the proposed method over the reconstruction-based approach, the comparison results of different AE structures and memory modules are shown in

TABLE V
THE RESULTS OF THE COMPARISON BETWEEN THE PROPOSED METHOD AND BASELINES. THE BEST RESULTS ARE MARKED IN BOLD. THE RESULTS OF 'AE+SKIP' ARE FROM [11]. 'W/O H' REPRESENTS THE PROPOSED METHOD WITHOUT HISTOGRAM ERROR ESTIMATION MODULE.

| Method | Bottle | Hazelnut | Capsule | Metal Nut | Leather | Pill | Wood | Carpet | Tile | Grid | Cable | Transistor | Toothbrush | Screw | Zipper | Mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| AE-SSIM [14] | 88 | 54 | 61 | 54 | 46 | 60 | 83 | 67 | 52 | 69 | 61 | 52 | 74 | 51 | 80 | 63 |
| AE-$L_2$ [14] | 80 | 88 | 62 | 73 | 44 | 62 | 74 | 50 | 77 | 78 | 56 | 71 | 98 | 69 | 80 | 71 |
| AE+skip | 71.3 | 82.8 | 74.7 | 33.6 | 57.0 | 85.3 | 97.7 | 38.5 | 98.6 | 87.9 | 57.9 | 74.9 | 74.2 | 100.0 | 69.6 | 73.6 |
| MemAE [10] | 95.4 | 89.1 | 83.1 | 53.7 | 61.1 | 88.3 | 95.4 | 45.4 | 63.0 | 94.6 | 69.4 | 79.3 | 97.2 | 99.2 | 87.1 | 80.2 |
| DAAD [11] | 97.5 | 89.3 | 86.6 | 55.2 | 62.8 | 89.8 | 95.7 | 67.1 | 82.5 | 97.5 | 72.0 | 81.4 | 98.9 | 100.0 | 90.6 | 84.5 |
| DAAD+ [11] | 97.6 | 92.1 | 76.7 | 75.8 | 86.2 | 90.0 | 98.2 | 86.6 | 88.2 | 95.7 | 84.4 | 87.6 | 99.2 | 98.7 | 85.9 | 89.5 |
| **OURS (W/O H)** | 98.7 | 94.8 | 72.5 | 75.3 | 97.2 | 69.3 | 100.0 | 95.0 | 88.1 | 99.1 | 81.6 | 90.8 | 87.5 | 100 | 82.43 | 88.8 |
| **OURS** | 95.2 | 99.4 | 82.3 | 84.5 | 94.5 | 86.1 | 100.0 | 93.1 | 97.2 | 97.1 | 85.6 | 92.4 | 95.8 | 97.0 | 77.3 | **91.8** |

TABLE VI
THE INFLUENCE OF DIFFERENT VALUES OF $\delta_H$ FOR ANOMALY DETECTION.

| $\delta_H$ | 0% | 25% | 50% | 75% | 90% |
|---|---|---|---|---|---|
| AUROC | 88.81 | 90.1 | 90.56 | 91.05 | 89.96 |

Table V. The different baselines are set up as follows :

- *AE-SSIM:* Vanilla AE structure using SSIM loss function.
- *AE-$L_2$:* Vanilla AE structure using $L_2$ loss function.
- *AE+skip:* Vanilla AE structure using $L_2$ loss function with skip connections.
- *MemAE:* Vanilla AE structure using $L_2$ loss function with traditional memory module (Figure 1 (1)).
- *DAAD:* Vanilla AE is equipped with multi-scale block-wise memory module (Figure 1 (2)).
- *DAAD+ :* DAAD+ is DAAD complemented with the adversarially learned representation.
- *OURS W/O H:* Vanilla AE is equipped with PMB module.
- *OURS:* Vanilla AE is equipped with PMB module and parameter-free Histogram Error Estimation module.

The experimental results show that the proposed method improves 27%-20% compared to the vanilla AE method, and improves compared to both MemAE and DAAD methods. These indicate the effectiveness of the novel query generation method of the PMB module, which allows the memory module to store normal features. As in the 'WOOD' class, AE+skip performance can reach 97.7%, while MemAE, DAAD, DAAD+ all reduce the performance, and the proposed method further improves the anomaly detection performance. Comparing the results of MemAE, DAAD, and OURS W/O H, these show that the proposed query (Figure 1 (3)) is the most effective among the three query generation methods shown in Figure 1.

**Impact of Histogram Error Estimation module.** To verify the effect of the proposed Histogram Error Estimation module for anomaly detection and localization, experiments are conducted by varying $\delta_H$ within [0%, 25%, 50%, 75%, 90%]. The anomaly detection results are shown in Table VI. When
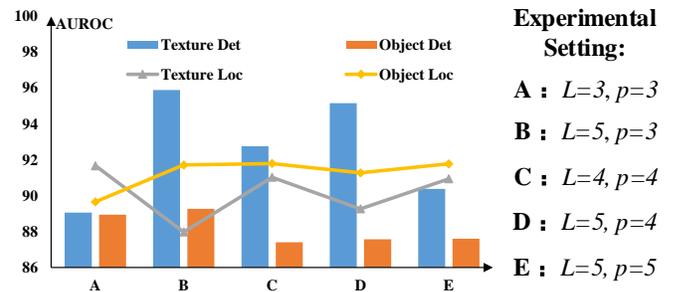


Fig. 9. The influences of the combination of skip connection and PMB on anomaly detection and anomaly localization. 'Det' denotes anomaly detection and 'Loc' denotes anomaly localization.

$\delta_H = 0\%$, which means that the Histogram Error Estimation module is not exploited, the detection performance dramatically decreases. As $\delta_H$ is increased to 75%, the performance of anomaly detection is improved and is 2% higher than the performance without the error estimation module. When $\delta_H = 90\%$, the result of anomaly detection decreases to 89.96%. The overall results in Table VI show that the Histogram Error Estimation module can well estimate the unfavorable reconstruction error caused by AE and effectively improve the performance of anomaly detection. However, when $\delta_H$ is close to 90%, the performance of anomaly detection is decreased. The reason is anomaly region errors affecting the estimation of $\mu$, resulting in a large estimation bias in $\mu$. For anomaly localization, after eliminating the reconstruction error caused by AE, $A'$ highlights the favorable reconstruction error of the anomaly regions caused by the PMB module, which significantly improves the anomaly localization. From Figure 8, it can be found that the localization performance of all categories is greatly improved when the proposed histogram error estimation module is introduced. By comprehensively considering the performance of anomaly localization and detection, $\delta_H$ is set to 50% in experiments.

**Impacts of $L$ and $p$.** In the proposed method, the combination strategy of the proposed PMB module and the skip connection are utilized. Experiments are conducted on the MVTec AD

TABLE VII
THE INFLUENCE OF DIFFERENT EXPERIMENTAL SETTINGS OF $k$ VALUES
FOR MVTEC AD.

| MVTec AD | k0 | k1 | k2 | k3 | k4 | k5 | k6 |
|---|---|---|---|---|---|---|---|
| Det | 85.6 | 83.9 | 89.6 | 90.7 | 88.8 | 90.4 | 91.5 |
| Loc | 93.0 | 91.1 | 91.1 | 90.8 | 89.7 | 89.1 | 89.5 |

TABLE VIII
THE INFLUENCE OF DIFFERENT $k$ EXPERIMENTAL SETTINGS ON ANOMALY
DETECTION IN MNIST AND RETINAL-OCT DATASETS.

| Det | k=(1,1,1) | k=(2,2,2) | k=(4,4,4) | k=(8,8,1) |
|---|---|---|---|---|
| MNIST | 98.1 | 96.5 | 97.6 | 95.8 |
| Retinal-OCT | 96.8 | 98.0 | 96.8 | - |

dataset to study the sensitiveness of $L$ and $p$. The results are shown in Figure 9. From the results by varying $p$ from 3 to 5 with $L = 5$, it can be seen that the alternate utilization of the skip connection and PMB can improve the detection performance, especially for texture images. Compared with the setting 'E' which does not use the skip connection, the performance by using settings 'B' and 'D' is improved by 6%. Because the detailed features provided by the PMB module guarantee the reconstruction ability of samples while the skip connection provides high-level information to keep the category consistent. For anomaly localization, it is found that the larger the $p$, that is, the more PMBs are introduced, the better the localization performance. In particular, the performance changes more significantly in the texture images, indicating that multiple PMBs can learn detailed normal features more effectively for subtle anomaly localization.

**Impact of $k$.** To evaluate the effect of the number of partitions $k$ on anomaly detection, experiments with different settings were conducted on three datasets. In the MVTec AD dataset, the partitions were set to k0: $(2, 2, 2, 2)$, k1: $(8, 8, 8, 8)$, k2: $(16, 16, 16, 8)$, k3: $(32, 16, 16, 8)$, k4: $(32, 32, 16, 8)$, k5: $(32, 16, 16, 16)$, k6: $(32, 32, 16, 16)$. $\delta_H$ is taken as 50%. The results of anomaly detection and localization on MVTec AD are shown in Table VII. The results of anomaly detection on the remaining datasets are shown in Table VIII. For the anomaly detection task, on the MVTec AD dataset, The larger $k$ the better the anomaly detection performance. Conversely, the smaller $k$ for MNIST and Retinal-OCT datasets, the better the detection performance. The reason is that the MVTec dataset has a large resolution and carries a large amount of information. Thus multiple storage units are adapted to store the detailed features to ensure the quality of normal samples reconstruction. While MNIST and Retinal-OCT have low resolution and simple semantic information, too large $k$ will make the anomaly features easy to generalize. Therefore, a smaller $k$ is adopted. For the anomaly localization task, the experimental results show that the smaller the $k$ value, the better the localization performance, while the larger the $k$ value, the smaller the impact on localization. The reason is that the expressiveness is relatively weaker when $k$ becomes

smaller, and the anomalous features are less likely to be expressed as normal features. Therefore, after eliminating the errors with the histogram error estimation module, the anomaly localization performance is excellent. It is worth noting that the values of $k$ for anomaly localization and anomaly detection are not contradictory. Larger $k$ values indicate a larger number of partitions, more detailed features can be learned, and the overall reconstruction error of normal samples is smaller. Therefore, a large $k$ value is more suitable for anomaly detection. A smaller $k$ value learns limited normal features and better amplifies the error after reconstruction of normal and abnormal pixels. Due to the Histogram Error Estimation module, excellent localization results can still be achieved after adaptive elimination of the unfavorable error.

## V. CONCLUSION

This paper proposes a new unsupervised visual anomaly detection method by jointly exploring AE and a novel memory module. To address the problem of the existing memory module, a new PMB module with a novel query generation method can learn and store the features of normal samples. With the successful reconstruction of normal features, it makes abnormal features only fit the normal features stored in PMB, which results in the reconstruction error of abnormal regions being larger. Furthermore, to eliminate the cumulative reconstruction error caused by AE, a novel Histogram Error Estimation module is proposed by exploring the reconstruction errors of normal regions. The detection and localization performance is significantly improved. Finally, we explored the optimal combination of memory module and skip connections. Extensive experiments are conducted on three widely-used datasets to verify the effectiveness of the proposed method for visual anomaly detection. In the future, we will explore a more suitable model to estimate the error caused by AE and will explore the generalizability of the proposed PMB module to other computer vision tasks.

## REFERENCES

[1] I. Golan and R. El-Yaniv, "Deep anomaly detection using geometric transformations," in *Proceedings of the Annual Conference on Neural Information Processing Systems*, 2018, pp. 9781–9791.

[2] A. Zimek, E. Schubert, and H.-P. Kriegel, "A survey on unsupervised outlier detection in high-dimensional numerical data," *Statistical Analysis and Data Mining: The ASA Data Science Journal*, vol. 5, no. 5, pp. 363–387, 2012.

[3] K. Leung and C. Leckie, "Unsupervised anomaly detection in network intrusion detection using clusters," in *Proceedings of the Australasian conference on Computer Science*, 2005, pp. 333–342.

[4] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger, "Mvtec ad–a comprehensive real-world dataset for unsupervised anomaly detection," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9592–9600.

[5] X. Mo, V. Monga, R. Bala, and Z. Fan, "Adaptive sparse representations for video anomaly detection," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 24, no. 4, pp. 631–645, 2014.

[6] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger, "Uninformed students: Student-teacher anomaly detection with discriminative latent embeddings," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2020, pp. 4183–4192.

[7] S. Mei, Y. Wang, and G. Wen, "Automatic fabric defect detection with a multi-scale convolutional denoising autoencoder network model," *Sensors*, vol. 18, no. 4, p. 1064, 2018.

[8] S. Zhang, M. Gong, Y. Xie, A. K. Qin, H. Li, Y. Gao, and Y.-S. Ong, "Influence-aware attention networks for anomaly detection in surveillance videos," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 8, pp. 5427–5437, 2022.

[9] Z. Li, C. Wang, M. Han, Y. Xue, W. Wei, L.-J. Li, and L. Fei-Fei, "Thoracic disease identification and localization with limited supervision," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8290–8299.

[10] 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 *Proceedings of the International Conference on Computer Vision*, 2019, pp. 1705–1714.

[11] J. Hou, Y. Zhang, Q. Zhong, D. Xie, S. Pu, and H. Zhou, "Divide-and-assemble: Learning block-wise memory for unsupervised anomaly detection," pp. 8771–8780, 2021.

[12] S. Zhai, Y. Cheng, W. Lu, and Z. Zhang, "Deep structured energy based models for anomaly detection," in *Proceedings of the International Conference on Machine Learning*. PMLR, 2016, pp. 1100–1109.

[13] B. Zong, Q. Song, M. R. Min, W. Cheng, C. Lumezanu, D. Cho, and H. Chen, "Deep autoencoding gaussian mixture model for unsupervised anomaly detection," in *Proceedings of the International Conference on Learning Representations*, 2018.

[14] P. Bergmann, S. Löwe, M. Fauser, D. Sattlegger, and C. Steger, "Improving unsupervised defect segmentation by applying structural similarity to autoencoders," in *Proceedings of the International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications*, 2019, pp. 372–380.

[15] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.

[16] S. Wold, K. Esbensen, and P. Geladi, "Principal component analysis," *Chemometrics and intelligent laboratory systems*, vol. 2, no. 1-3, pp. 37–52, 1987.

[17] L. Xiong, B. Póczos, and J. Schneider, "Group anomaly detection using flexible genre models," 2011.

[18] M. Salehi, A. Arya, B. Pajoum, M. Otoofi, A. Shaeiri, M. H. Rohban, and H. R. Rabiee, "Arae: Adversarially robust training of autoencoders improves novelty detection," *arXiv preprint arXiv:2003.05669*, 2020.

[19] H. Park, J. Noh, and B. Ham, "Learning memory-guided normality for anomaly detection," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14372–14381.

[20] A. Kumar, O. Irsoy, P. Ondruska, M. Iyyer, J. Bradbury, I. Gulrajani, V. Zhong, R. Paulus, and R. Socher, "Ask me anything: Dynamic memory networks for natural language processing," in *Proceedings of the International Conference on Machine Learning*. PMLR, 2016, pp. 1378–1387.

[21] C. Fan, X. Zhang, S. Zhang, W. Wang, C. Zhang, and H. Huang, "Heterogeneous memory enhanced multimodal attention model for video question answering," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1999–2007.

[22] W. Lu, J. Zhang, X. Zhao, W. Zhang, and J. Huang, "Secure robust JPEG steganography based on autoencoder with adaptive BCH encoding," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 7, pp. 2909–2922, 2021.

[23] Z. Liu, W. Siu, and Y. Chan, "Photo-realistic image super-resolution via variational autoencoders," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 4, pp. 1351–1365, 2021.

[24] Y. Fei, C. Huang, C. Jinkun, M. Li, Y. Zhang, and C. Lu, "Attribute restoration framework for anomaly detection," *IEEE Transactions on Multimedia*, 2020.

[25] T. Wang, X. Xu, F. Shen, and Y. Yang, "A cognitive memory-augmented network for visual anomaly detection," *Journal of Automatica Sinica*, vol. 8, no. 7, pp. 1296–1307, 2021.

[26] B. Schölkopf, A. J. Smola, F. Bach *et al.*, *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press, 2002.

[27] Y. Chen, X. S. Zhou, and T. Huang, "One-class svm for learning in image retrieval," in *Proceedings of the International Conference on Image Processing (Cat. No.01CH37205)*, vol. 1, 2001, pp. 34–37 vol.1.

[28] L. Ruff, R. Vandermeulen, N. Goernitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, and M. Kloft, "Deep one-class classification," in *Proceedings of the International Conference on Machine Learning*. PMLR, 2018, pp. 4393–4402.

[29] R. Chalapathy, A. K. Menon, and S. Chawla, "Anomaly detection using one-class neural networks," *arXiv preprint arXiv:1802.06360*, 2018.

[30] J. Kim and K. Grauman, "Observe locally, infer globally: a space-time mrf for detecting abnormal activities with incremental updates,"

[31] W. Li, V. Mahadevan, and N. Vasconcelos, "Anomaly detection and localization in crowded scenes," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 1, pp. 18–32, 2014.

[32] M. Salehi, N. Sadjadi, S. Baselizadeh, M. H. Rohban, and H. R. Rabiee, "Multiresolution knowledge distillation for anomaly detection," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14902–14912.

[33] G. Wang, S. Han, E. Ding, and D. Huang, "Student-teacher feature pyramid matching for unsupervised anomaly detection," *arXiv preprint arXiv:2103.04257*, 2021.

[34] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," *Advances in neural information processing systems*, vol. 27, 2014.

[35] 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 *Proceedings of the International conference on information processing in medical imaging*. Springer, 2017, pp. 146–157.

[36] S. Akcay, A. Atapour-Abarghouei, and T. P. Breckon, "Ganomaly: Semi-supervised anomaly detection via adversarial training," in *Proceedings of the Asian Conference on Computer Vision*. Springer, 2018, pp. 622–637.

[37] Y. Zhang, X. Nie, R. He, M. Chen, and Y. Yin, "Normality learning in multispace for video anomaly detection," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 9, pp. 3694–3706, 2021.

[38] M. Hasan, J. Choi, J. Neumann, A. K. Roy-Chowdhury, and L. S. Davis, "Learning temporal regularity in video sequences," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2016, pp. 733–742.

[39] W. Liu, W. Luo, D. Lian, and S. Gao, "Future frame prediction for anomaly detection–a new baseline," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6536–6545.

[40] W. Luo, W. Liu, and S. Gao, "A revisit of sparse coding based anomaly detection in stacked rnn framework," in *Proceedings of the International Conference on Computer Vision*, 2017, pp. 341–349.

[41] A.-S. Collin and C. De Vleeschouwer, "Improved anomaly detection by training an autoencoder with skip connections on images corrupted with stain-shaped noise," in *Proceedings of the International Conference on Pattern Recognition*. IEEE, 2021, pp. 7915–7922.

[42] J. Weston, S. Chopra, and A. Bordes, "Memory networks," *arXiv preprint arXiv:1410.3916*, 2014.

[43] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[44] X. Wen, Z. Han, and Y. Liu, "CMPD: using cross memory network with pair discrimination for image-text retrieval," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 6, pp. 2427–2437, 2021.

[45] 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.

[46] N. Cohen and Y. Hoshen, "Sub-image anomaly detection with deep pyramid correspondences," *CoRR*, vol. abs/2005.02357, 2020.

[47] D. Dehaene and P. Eline, "Anomaly localization by modeling perceptual features," *CoRR*, vol. abs/2008.05369, 2020. [Online]. Available: https://arxiv.org/abs/2008.05369

[48] C.-L. Li, K. Sohn, J. Yoon, and T. Pfister, "Cutpaste: Self-supervised learning for anomaly detection and localization," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2021, pp. 9664–9674.

[49] S. Venkataramanan, K.-C. Peng, R. V. Singh, and A. Mahalanobis, "Attention guided anomaly localization in images," in *Proceedings of the European Conference on Computer Vision*. Springer, 2020, pp. 485–503.

[50] M. Salehi, A. Eftekhar, N. Sadjadi, M. H. Rohban, and H. R. Rabiee, "Puzzle-ae: Novelty detection in images through solving puzzles," *CoRR*, vol. abs/2008.12959, 2020.

[51] A. Collin and C. D. Vleeschouwer, "Improved anomaly detection by training an autoencoder with skip connections on images corrupted with stain-shaped noise," in *Proceedings of the International Conference on Pattern Recognition,*, 2020, pp. 7915–7922.

[52] X. Li, I. Kiringa, T. Yeap, X. Zhu, and Y. Li, "Exploring deep anomaly detection methods based on capsule net," in *Proceedings of the Canadian Conference on Artificial Intelligence*. Springer, 2020, pp. 375–387.

in *Proceedings of the conference on Computer Vision and Pattern Recognition*. IEEE, 2009, pp. 2921–2928.

[53] P. Perera, R. Nallapati, and B. Xiang, "Ocgan: One-class novelty detection using gans with constrained latent representations," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2898–2906.

[54] D. Abati, A. Porrello, S. Calderara, and R. Cucchiara, "Latent space autoregression for novelty detection," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2019, pp. 481–490.

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

[56] Y. LeCun, "The mnist database of handwritten digits," *http://yann. lecun. com/exdb/mnist/*, 1998.

[57] J. Yang, K. Zhou, Y. Li, and Z. Liu, "Generalized out-of-distribution detection: A survey," *CoRR*, vol. abs/2110.11334, 2021.

[58] P. Gholami, P. Roy, M. K. Parthasarathy, and V. Lakshminarayanan, "Octid: Optical coherence tomography image database," *Computers & Electrical Engineering*, vol. 81, p. 106532, 2020.

[59] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Proceeding of the International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

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

[61] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the international Conference on Computer Vision*, 2017, pp. 2223–2232.

[62] K. Zhou, Y. Xiao, J. Yang, J. Cheng, W. Liu, W. Luo, Z. Gu, J. Liu, and S. Gao, "Encoding structure-texture relation with p-net for anomaly detection in retinal images," in *Proceeding of the European Conference of Computer Vision*. Springer, 2020, pp. 360–377.

[63] Y. Liu, C. Zhuang, and F. Lu, "Unsupervised two-stage anomaly detection," *arXiv preprint arXiv:2103.11671*, 2021.

[64] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.

[65] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[66] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.

[67] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *Proceedings of the International Conference on Machine Learning*. PMLR, 2019, pp. 6105–6114.

**Zechao Li** is currently a Professor at the Nanjing University of Science and Technology. He received his Ph.D degree from National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences in 2013, and his B.E. degree from the University of Science and Technology of China in 2008. His research interests include big media analysis, computer vision, etc. He was a recipient of the best paper award in ACM Multimedia Asia 2020, and the best student paper award in ICIMCS 2018. He serves as an Associate Editor for IEEE TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS.

**Peng Xing** is currently pursuing the master's degree with the School of Computer Science and Engineering, Nanjing University of Science and Technology. His current research interests include anomaly detection and unsupervised deep learning.